## Composite Endpoints

Dean Follmann National Institute of Allergy and Infectious Diseases

### Outline

- Simple Phase III Trial
- Composite endpoints
- Co-equal primary endpoints
- Co-equal surrogate endpoints
- Combining 1<sup>o</sup> & 2<sup>o</sup> endpoints
- Ranking Methods
- Conclusions

#### Simple Phase III trial

- Choose a single relevant endpoint
  - Death
  - Ejection fraction of the left ventricle
- Size trial for 90% power to detect a clinically important effect
  - 20% reduction in mortality
  - .05 difference in EF

#### Sample size formula

- Treatment Effect: more is better
  - Pick a responsive endpoint
- Variability: less is better
  - Get more measurements/stable endpoint
- Events: more are better
  - Include sicker patients
  - Lengthen follow-up

Characteristics of a clinical endpoint (Neaton et al 1994)

- Should be relevant and easy to interpret.
- Should be clinically apparent and easy to diagnose.
- Should be sensitive to treatment differences.

#### More complicated world

- Occasionally, a single primary endpoint undesirable. Why?
  - Clinically important events are rare.
  - Effect of treatment manifested on a variety of important endpoints.

#### Example: ACES trial

- ACES—trial to evaluate antibiotics versus placebo in patients at risk of CHD events.
- Primary endpoint is
  - Hospitalization for unstable angina
  - CHD death
  - Nonfatal MI
  - Revascularization

#### **Composite Concerns**

- With a composite endpoint, relative importance of various constituent endpoints determined by frequency.
- CHD death or revascularization
  - CHD death 1%
  - Revascularization 10%

#### **Composite Concerns**

- Only include constituent endpoints who are reasonably influenced by treatment.
  - Treatment: 50% on death, 20% on MI
    - Control rate Treatment rate
  - Death .01 .005
  - MI .01 .008
- Death alone vs Death or MI: same power

#### Bonferroni approach.

- Use p-values for the two endpoints.
- Reject if p<sub>1</sub> or p<sub>2</sub> less than .05/2
- Inference drawn for each endpoint
- Good if treatment has entire effect on one endpoint or the other, don't know which one.

#### Example: PEPI

- Postmenopausal Estrogen/Progestin Interventions Trial. HRT's effect on risk factors for heart disease.
- 875 women assigned to 5 combinations.
- Primary endpoints
  - HDL-C
  - SBP
  - Serum insulin
  - fibrinogen

### O'Brien (1984) Rank-Sum method

- Rank each outcome and calculate an average rank for each patient
- See if average rank differs between groups.

Sub	X <sub>1</sub>	R <sub>1</sub>	X <sub>2</sub>	$R_2$	Avg R
Fred	3.3	2	87	1	1.5
Joe	4.1	3	105	2	2.5
Sam	1.7	1	1000	3	2.0

#### O'Brien OLS method

 Standardize each endpoint. Compute the average endpoint for each person and perform a t-test on the averages.

Sub	X <sub>1</sub>	(X <sub>1</sub> -O)/•	X <sub>2</sub>	(X <sub>2</sub> -O)/•	avg	
Fred	3.3	.22	87	72	25	
Joe	4.1	.87	105	42	.23	
Sam	1.7	-1.10	200	1.14	.03	

#### O'Brien GLS method

- Assume common treatment effect I solution
  - e.g. 1 standard deviation on both endpoints.
- Calculate a statistically optimal estimate of using a weighted average. (more correlated endpoints, less weight).
- Pocock Geller Tsiatis (1987) generalize to binary/survival etc endpoints.
- Many other methods conceptually similar: specify a model with the same for many endpoints.

#### Latent Variable models

- Assume each person has an underlying severity, S, which influences several endpoints.
- E.g. MPS---Lysosomal enzyme deficency
  - FVC
  - 6 minute walk
  - AHI
  - shoulder flexion
  - visual acuity
- Test whether the distribution of underlying severities is moved by treatment.

Conceptual framework for latent variable model



# Simple Model Y\_{i1} = B\_01 + D Z\_i + b\_i + e\_{i1} Y\_{i2} = B\_02 + D Z\_i + b\_i + e\_{i2} e\_ij ~ N(0, Vej) S\_i ~ N(0, Vs)

A model



- Hotelling T<sup>2</sup>---multivariate t-test
- Good for *any* treatment effect, so less good for uniformly beneficial treatment effects.

Rejection Region for Hotelling's T<sup>2</sup> Test



#### Rejection Regions for Hotelling's T<sup>2</sup> Test & O'Brien test



Rejection Regions for Hotelling's T<sup>2</sup>, O'Brien, & Bonferroni Tests



## Tilley et al (1996) for Stroke trial

- Trial of t-PA versus placebo in patients with acute ischemic stroke.
- Dichotomized 4 stroke scales.
- Discussed use of Bonferroni, Hotelling's Test & O'Brien's GLS test.
- Reject if

■ Mean(Z) > 1.96 \* [ (1+ 3 □ ) < ] </p>

Combining co-equal but surrogate endpoints

- Suppose both endpoints are surrogates.
- Ideally form a risk score.
  - $R = w_1 DBP + w_2 SBP + w_3 serum insulin...$
  - $R = w_1$  Hepatitis +  $w_2$  sex for drugs + ...

Do a t-test using R.

### Combining 1<sup>o</sup> and auxiliary...

- 1<sup>o</sup> endpoint alone: use Wilcoxon Rank sum approach.
- Compare each pair of treatment/control patients
- Image: marginal statement = 1 if "i" (in T) lives past "j" (in P)
- $Y_{ij} = 1/2$  if both live
- if "i" (in T) dies before "j" (in P)
- Form mean(Y<sub>ij</sub>) = Pr(live longer on T than P)
- Equivalent to ranking by death time.

#### Combining 1<sup>o</sup> and auxiliary...

- If both live, replace ½ with
  - p<sub>ij</sub> = Pr( i lives longer than j | CD4s)
- May be useful if
  - CD4/death relationship in past = future
  - Treatment effects CD4 counts & they differ at end
- Similar approach taken by Faucett Schenker Taylor (2002) who imputed death times.

#### "Utility" Ranking Methods

- May be hard to say MI is half as bad as death. But clearly death is worse.
  - Death is worst
    - Rank by death time
  - 2 Strokes worse than 1
    - Rank by time of first stroke
  - I Stroke worse than nothing
    - Rank by time of stroke.
- Compare the ranks between groups



1 * @	D S H S H	2 *D S H H	S * 1 * 1	He He	4	6 8 27 29 35 16 28 1 40 13 39 24 21	Mean 1-2 1-9 3-0 4-0 5-3 6-9 7-5 8-6 8-8 8-8 11-5 11-7 12-2 12-2 12-2	Std 04 05 03 09 15 14 18 16 18 18 27 31	death 04 04 08 1·1 2·0 2·1 2·9 3·0 2·8 3·0 3·3 2·8	events 1 1 1 1 1 1 2 3 3 2 2 2 2 1
* @ 	D S S H S H	* S H H	S * 1 * 1 *	р е не		6 8 27 29 35 16 28 1 40 13 39 24 21	1-2 1-9 3-0 4-0 5-3 6-9 7-5 8-6 8-8 11-5 11-7 12-2 12-2	04 05 03 09 15 14 18 16 18 18 18 27 31	04 04 08 1·1 2·0 2·1 2·9 3·0 2·8 3·0 3·3 2·8	1 1 1 2 3 2 2 2 1
œ H	D S H S H	* D S H H	S * 1 * 1	р е не		8 27 29 35 16 28 1 40 13 39 24 21	1-9 3-0 4-0 5-3 6-9 7-5 8-6 8-8 11-5 11-7 12-2 12-2	05 03 09 15 14 18 16 18 18 27 31	04 08 1·1 2·0 2·1 2·9 3·0 2·8 3·0 3·3 3·3 2·8	1 1 2 3 3 2 2 2 1
ee H	D S H S H	т В Н Н	S * 1 * 1 *	) е Не		27 29 35 16 28 1 40 13 39 24 21	3-0 4-0 5-3 6-9 7-5 8-6 8-8 11-5 11-7 12-2 12-2	03 09 15 14 18 16 18 18 18 27 31	0-8 1-1 2-0 2-1 2-9 3-0 2-8 3-0 3-3 3-3 2-8	1 1 2 3 2 2 2 1
. <b>H</b>	ы s н s H	* S Н Н	S * 1 * 1 *	) е Не		29 35 16 28 1 40 13 39 24 21	40 5-3 6-9 7-5 8-6 8-8 11-5 11-7 12-2 12-2	0-3 0-9 1-5 1-4 1-8 1-6 1-8 1-8 1-8 2-7 3-1	1·1 2·0 2·1 2·9 3·0 2·8 3·0 3·3 3·3 2·8	1 2 3 3 2 2 2 1
. <b>H</b>	s н s	т s н н	s * 1 * 1	) @ Н@	•	35 16 28 1 40 13 39 24 21	5-3 6-9 7-5 8-6 8-8 11-5 11-7 12-2 12-2	09 1.5 1.4 1.8 1.6 1.8 1.8 2.7 3.1	20 21 29 30 28 30 33 33	1 2 3 3 2 2 2 1
Н	з н s H	ы н н	s * * 1	) е не	•	16 28 1 40 13 39 24 21	6-9 7-5 8-6 8-8 11-5 11-7 12-2 12-2	1.5 1.4 1.8 1.6 1.8 1.8 2.7 3.1	2-1 2-9 3-0 2-8 3-0 3-3 2-8	2 3 2 2 2 1
H	H S H	s н	* * 1 * 1 *	р е не и	•	28 - 1 40 13 39 24 21	7.5 8.6 8.8 11.5 11.7 12.2 12.2	1.4 1.8 1.6 1.8 1.8 2.7 3.1	2-9 3-0 2-8 3-0 3-3 2-8	3 2 2 2 1
H	н s н	H H	•	с С С С С С С С С С С С С С С С С С С С	•	1 40 13 39 24 21	8.8 11.5 11.7 12.2 12.2	1.8 1.6 1.8 1.8 2.7 3.1	3-0 2-8 3-0 3-3 2-8	2 2 2 1
ur kriti Kritik	н s н	H	•	е не		40 13 39 24 21	8.8 11.5 11.7 12.2 12.2	1.8 1.8 2.7 3.1	2.8 3-0 3-3 2.8	2 2 2 1
	s н	н	٠	е не це		13 39 . 24 21	11-5 11-7 12-2 12-2	1.8 1.8 2.7 3.1	3-0 3-3 2-8	2
	s н	•	•	e He		39 24 21	11-7 12-2 12-2	1·8 2·7 3·1	3.3	2
	в		•	Н@ @		24	12.2	2.7	2.8	1
	к н			на 		21	12-2	3.1		
	н			uu@					3.7	3
						2	12.5	1.0	3-5	2
				ппе		38	15-2	2.3	3.4	3
				D		10	16-7	2.1	3.5	1
				@		36	17.6	2.5	3.6	. 1
	-		-	*		34	18.5	1.6	3.7	1
	, S		S	S		32	21-0	8.1		3
S	_	н				15	21.5	1.9		2
н	S					42	21.7	1.9		. 2
н	н					3	23.5	4.7		2
нн				-		19	23.8	4-6		2
S				S		4	24.2	40		2
	нн	н				14	24.7	4.4		3
	SS					20	25.2	4.6		2
	н	S				23	27.7	3-7		2
S		_	-			18	28.4	5-3		1
	1	s	S			17	28-6	4.4		2
	н			н	_	26	30-0	4.4		2
		н		_	S	30	30-9	3.5		2
			н	S		41	31.9	3.7		2
			SH	I		9	32.2	3-6		2
н				_		12	32.4	4.7		1
				SS		37	34-0	3.5		2
	н					33	35-1	3-4		1
			н	H		31	35.7	2.8		2
			S			7	35.8	3-5		1
				S		22	38.2	2.7		1
				н		43 -	39-5	2.9		1
				S		5	39.8	20		1
				н		11	41.4	2.5		1
						25	450	0-0		0
	н	н	H H	H H S	H SS H H H S H S H	H H SS H S H S H H L 2 3 4	H 12 H SS 37 H H 31 S 7 H 43- S 5 H 11 25	H 12 32.4 H 12 32.4 S H 9 32.2 33 35-1 H 33 35-1 H 13 557 S 7 35-8 S 7 35-8 S 22 38-5 H 43 - 39-5 S 5 39-8 H 11 41-4 25 43-0 1 2 3 4	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$

Table III. The cards ordered by consensus ranks and the mean and standard deviation of the 19 ranks

The mean and std (standard deviation) of the rank are based on the sample of 19 rankers.

- H non-fatal heart attack
- @ fatal heart attack
- # fatal stroke
- S non-fatal stroke
- D death from neither heart attack nor stroke

#### HIV vaccine trials

- Want HIV vaccine to reduce acquistion and also post-infection viral load for those infected. How to combine?
  - Those who are uninfected get best rank
  - Those who are infected are ranked by viral load "setpoint" lower setpoints get higher ranks.

## Weighting

- You may not be interested in weighting, but weighting is interested in you.
- Approaches we discussed.
  - Equal weight for all endpoints (e.g. OLS)
  - More weight for frequent events (e.g. composite)
  - Less correlated outcomes more weight (e.g. GLS)
- Clinically interpretable weights?

#### Conclusions

- Common approaches are to pick a composite endpoint or adopt a Bonferroni correction.
- Clinical relevance / interpretability paramount.
- Appropriate approach depends heavily on the application.
- Novel endpoints/analysis approaches should be thoroughly investigated.

- Bandeen-Roche K, Miglioretti DL, Zeger SL, et al. (1997) "Latent variable regression for multiple discrete outcomes" Journal of the American Statistical Association, 1375-1386.
- Bjorling L, Hodges J, (1997) "Rule-Based Ranking schemes for antiretroviral trials" Statistics in Medicine, 1175-1191.
- Faucett C, Schenker N, Taylor J, (2002) "Survival Analysis using auxiliary variables via multiple imputation, with application to AIDS clinical trial data" Biometrics, 37-47.
- Follmann D, (1995) "Multivariate Tests for Multiple Endpoints in Clinical Trials" Statistics in Medicine, 1163-1176.
- Follmann D, (1996) "A Simple Multivariate Test for One Sided Alternatives" Journal of the American Statistical Association, 854-861.



- Lefkopoulou M, Ryan L, (1993) "Global tests for multiple binary outcomes" Biometrics, 975-988.
- Legler JM, Ryan LM, (1997) "Latent variable models for teratogenesis using multiple binary outcomes" Journal of the American Statistical Association, 13-20.
- Miller VT, Larosa J, Barnabei V, et al. (1995) "Effects of Estrogen or Estrogen/Progestin Regimens on Heart-Disease Risk Factors in Postmenopausal Women" The Postmenopausal Estrogen/Progestin Interventions (PEPI) Trial, 199-208.
- Neaton J, Wentworth D, Rhame, et al. (1994) "Considerations in choice of a clinical endpoint for AIDS clinical trials" Statistics in Medicine, 2107-2125.



- O'Brien PC, Geller NL, (1997) "Interpreting tests for efficacy in clinical trials with multiple endpoints" Controlled Clinical Trials, 222-227.
- Pocock S, Geller N, Tsiatis A, (1987) "The analysis of multiple endpoints in clinical trials" Biometrics, 487-498.
- Sammel M, Lin X, Ryan L, (1999) "Multivariate linear mixed models for multiple outcomes" Statistics in Medicine, 2479-2492.
- Sammel MD, Ryan LM, Louise M, Legler JM, (1997) "Latent variable models for mixed discrete and continuous outcomes" Journal of the Royal Statistical Society, Series B, Methodological, 667-678

- Tang DI, Geller NL, Pocock SJ, (1993) "On the design and analysis of randomized clinical trials with multiple endpoints" Biometrics, 23-30.
- Tilley BC, Marler J, Geller NL, et al. (1996) "Use of a global test for multiple outcomes in stroke trials with application to the national institute of neurological disorders and stroke t-PA stroke trial" Stroke, 2136-2142.
- Wassmer G, Reitmeir P, Kieser M, Lehmacher W, (1999) "Procedures for testing multiple endpoints in clinical trials: An overview" Journal of Statistical Planning and Inference, 69-81.
- Zhang J, Quan H, Ng J, et al. (1997) "Some statistical methods for multiple endpoints in clinical trials" Controlled Clinical Trials, 204-221.
- Zhong A, Song C, Reiss TF, (2004) "An endpoint for worsening asthma: Development of a sensitive measure and its properties" Drug Information Journal 5-13

#### Mariamman: Goddess of pox



- Afflicted individuals provide offerings
- Follow them home
- Successful ring vaccination
- Smallpox eradicated

#### Novel Design Issues

- Would a crossover trial make sense?
  - Area under EDSS curve over time.
- Enroll patients during a remission?
- For a phase II study, could a placebo be ethically used for a short while?
- Could all patients receive drug at end of study?
- Can we cross-over at time of failure?

#### Aldurazyme trial in MPS

- MPS: lysosomal enzyme deficency, leads to GAG accumulation with multisystemic effects.
- Inclusion criteria:
  - Stand 6 minutes, walk > 5 meters
- weekly IV infusion for ½ year.
- N=45
- Endpoints: FVC, 6 minute walk, AHI, shoulder flexion, visual acuity.

#### Phase 3 Study: Post-Hoc Analysis Composite Endpoint

#### Placebo

#### Aldurazyme

k - Krone

Not Available 61

Pate at	PEC	SMWT	SHFLEX	AHI	ACUITY		Patient	FFC	SHWT	SHFLET	j.
	11%	54m	20 deg	10 ev/br	2-lines			11%	5.4m	20 deg	1.0
1							24				
2							2.5				
- 3							2.6				
4							27				
6							28				
6							2.9				
т							- 20				
0							31				
9							32				
5.0							3.3				
11							24				
1.2							3.5				Γ
1.3							36				
1.0							37				
15							39				
16							39				
1.2							4.0				
5.0							41				
1.0							42				
2.0							4.9				
2.1							- 44				
2.2							45				
2.3											

Clinically Significant Changes



Improvement

Decline



#### Phase 3 Study: Post-Hoc Analysis Aldurazyme Leads to Net Improvement



#### Phase 3 Study: Aldurazyme Reduces Urinary GAG Levels



#### Example: asthma score

- Asthma: manifold symptoms, periodic worsening.
- Zhang, Song, Reiss (2004) proposed
  - PEF decrease >20%
  - 2+ puffs/day of beta-agnoist
  - Increase in symptom score > 50%
  - 3+ nighttime awakenings
  - PEF < 180 L/min</p>
  - Hospital visit
- Showed good correlation with other global evaluations.

#### Two endpoints—setup

- Let X<sub>1</sub> and X<sub>2</sub> be two endpoints.
  - Two stroke scales, DBP & SBP, time to AIDS/Death & CD4.
- Let Z<sub>1</sub> and Z<sub>2</sub> be the associated standardized test statistics.
  - E.g. two tests of proportions, two t-tests, log-rank & t-test.
- Let  $p_1$  and  $p_2$  be the two p-values.
- Let's assume X<sub>1</sub> and X<sub>2</sub> are independent

#### Ranking generalization

- Compare each pair of treatment/control patients
- = 1 if "i" (in T) does better "j" (in P)
- $Y_{ij} = 1/2$  if same
- $= 0 \quad \text{if "i" (in T) does worse "j" (in P)}$
- Can compare "i" & "j" over common followup.



