

## BACKGROUND

- Assessment of differences in event rates is a common endeavor in the evaluation of efficacy and safety of new treatments in clinical trials (in particular in oncology).
- We investigate the performance of different hypothesis tests for an overall survival endpoint.
- Stratified analyses are desired and sometimes even required by regulators.
- We illustrate the necessity of non-zero variance estimates – especially in the presence of strong prognostic stratification effects.
- Focus: comparison of event rates via Kaplan-Meier estimates for a pre-specified time  $t_0$
- One-sided test for superiority at significance level  $\alpha=0.05$ .

### GOALS:

- Choose best **stratum weights**
- Choose best Kaplan-Meier **variance estimator**
- Compare performance with **Cox model** (stratified)

## SIMULATION STUDY

- 1 experimental arm and 1 control arm
- Per treatment group: 65 subjects divided into 3 strata
- Time of interest:  $t_0 = 100$  (e.g., days)
- Survival time and censoring time from exponential distributions
- Constant censoring intensity  $\lambda_{cens}=0.005$  (leads to 22%-38% of patients censored in scenario 2 below)
- Hazard rates of **active treatment group** ( $G=1$ ) and **control group** ( $G=2$ ) in stratum  $k$ :  $\lambda_{k,G}$  for  $G=1,2$
- Two underlying models:
  - proportional hazard rates** that satisfy the **Cox model (COX)**

$$\lambda_{k,2}(t) = c_p \lambda_{k,1}(t)$$
 for all  $t$  and with same hazard ratio  $c_p \geq 1$  for all  $k$ .
  - additive survival difference (ASD)** at time  $t_0$ :
 
$$S_{k,2}(t_0) = S_{k,1}(t_0) - c_A$$
 for same difference in survival  $c_A \geq 0$  for all  $k$ .
- 10,000 simulation runs
- Simulate proportional effects (hazard ratios):  $c_p = 2.0, 2.5, 3.0, 3.5$
- Simulate additive effects:  $c_A = 0.10, 0.15, 0.20, 0.25$
- 3 different allocations in table below but results only illustrated for scenario 2

Scenario	$S_{k,1}(t_0)$ for $k=1,2,3$	n per treatment group for $k=1,2,3$	Note
1	0.80, 0.50, 0.30	20, 30, 15	Base case
2	<b>0.95, 0.70, 0.50</b>	<b>45, 10, 10</b>	<b>Largest stratum with greatest <math>S_{k,1}(t_0)</math></b>
3	<b>0.95, 0.70, 0.50</b>	<b>10, 10, 45</b>	<b>Smallest stratum with greatest <math>S_{k,1}(t_0)</math></b>

## Z-TESTS FOR KAPLAN-MEIER RATES

Z-test for 3 strata at time  $t_0$ :

$$Z = \frac{\sum_{k=1}^3 w_k (\hat{S}_{k,1} - \hat{S}_{k,2})}{\left(\sum_{k=1}^3 w_k^2 (\hat{\sigma}_{k,1}^2 + \hat{\sigma}_{k,2}^2)\right)^{0.5}}$$

- $\hat{S}_{k,G}$ : Kaplan-Meier estimator in stratum  $k$  and treatment group  $G=1,2$  at time  $t_0$
- $\hat{\sigma}_{k,G}^2$ : variance estimate of  $\hat{S}_{k,G}$ ,  $w_k$ : stratum weights

### Common choice:

**Greenwood** formula:

$$\hat{\sigma}_{k,G}^2 = \sum_{t_i \leq t_0} \frac{d_i}{n_i(n_i - d_i)}$$

**Problem:** can become 0 if no events observed or if all subjects have an event  $\rightarrow Z$  is not well-defined. In the simulations,  $Z$  is evaluated if at least one stratum with non-zero variance is present. Strata with 0 variance are excluded from the analysis.

### VARIANCE ESTIMATORS

#### Suggestion:

**Borkowf's adjusted hybrid variance estimator** [1]:

$$w_*(1 - w_*) / (n - m_c)$$

for  $w_* = (1 - n^{-1})w + (2n)^{-1}$ ,  $m_c$ : # censored subjects at  $t_0$ ,  $w$ : truncated version of  $\hat{S}_{k,G}$

Assures non-zero variances for all  $t$

### STRATUM WEIGHTS

**Inverse variance (IV) weight:**

$$w_k = \frac{(\hat{\sigma}_{k,1}^2 + \hat{\sigma}_{k,2}^2)^{-1}}{\sum_{j=1}^4 (\hat{\sigma}_{j,1}^2 + \hat{\sigma}_{j,2}^2)^{-1}}$$

**Motivation:** "least-squares", i.e., minimal variance in  $Z$  [2]

**Problem:** not defined for strata with 0 variance.

**Mantel-Haenszel (MH) weight:**

$$w_k = \frac{n_k m_k / (n_k + m_k)}{\sum_j (n_j m_j / (n_j + m_j))}$$

for  $n_k$  and  $m_k$  number of treatment and control subjects in stratum  $k$

**Motivation:** robust for sparse events [3] if  $n_k = m_k$ :  $w_k =$  subjects in stratum  $k$  / total number of subjects

2 underlying models with 2 different types of tests each:

	Test for ASD ( $H_0: c_A = 0$ )	Test for COX ( $H_0: c_p = 1$ )
Simulate: COX $c_p \geq 1$	<b>Z-tests assumptions violated X</b>	<b>Cox regression correct model ✓</b>
Simulate: ASD $c_A \geq 0$	<b>Z-tests correct model ✓</b>	<b>Cox regression assumptions violated X</b>

## RESULTS

- The **violation of assumptions** for the respective mis-specified stratified tests increases with increasing effect size in the underlying model: The effects  $c_p$  and  $c_A$  are not the same across all strata if the data are generated from the respective other model.

$c_p$	$c_p$ per stratum		
	stratum 1	stratum 2	stratum 3
0.10	3.168	1.432	1.322
0.15	4.350	1.676	1.515
0.20	<b>5.609</b>	<b>1.943</b>	<b>1.737</b>
0.25	6.954	2.239	2.000

**Table 1:** Proportional hazard ratios  $c_p$  in data from ASD models with actual effect  $c_A$ .

$c_A$	$c_A$ per stratum		
	stratum 1	stratum 2	stratum 3
2.0	0.048	0.210	0.250
2.5	0.070	0.290	0.323
3.0	<b>0.093</b>	<b>0.357</b>	<b>0.375</b>
3.5	0.114	0.413	0.412

**Table 2:** Additive survival effects  $c_A$  in data from COX models with actual effect  $c_p$ .

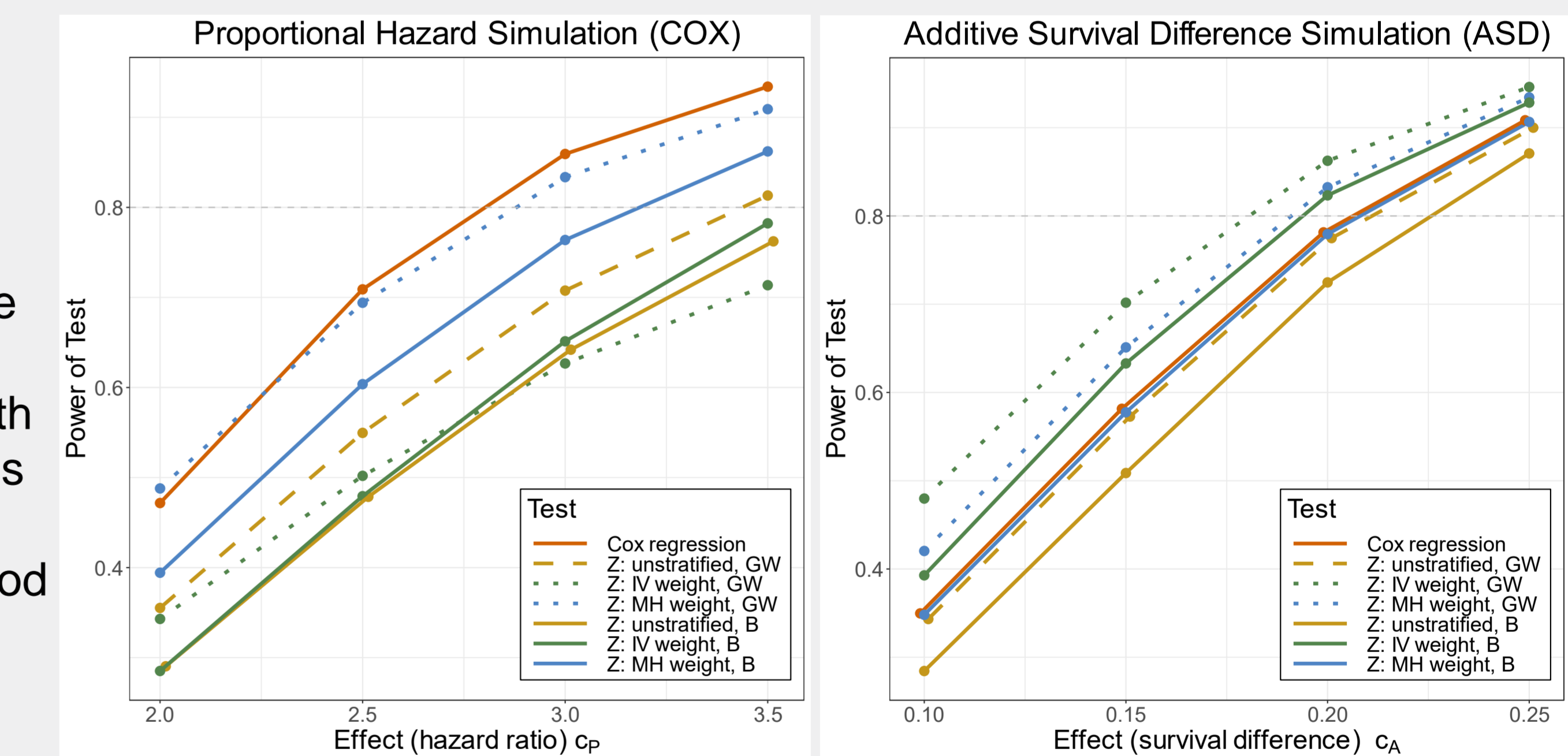
- At least one stratum with Greenwood variance equal to 0 occurred only in  $<0.01\%$  of all ASD simulation runs with  $c_A=0.2$  and in  $0.1\%$  of all COX simulations with  $c_p=3.0$  (in  $2.9\%$  of all simulations under  $H_0$ ). The case with 0 variance in all three strata did not occur in this scenario.
- In scenario 3 (not discussed here), up to  $45\%$  of all simulation runs had a stratum with 0 variance under  $H_0$ .

### Type I error:

- Cox regression controls type I error as expected.
- Stratified Z-tests with Greenwood variance inflate type I error and are excluded from the analysis of power below (dotted lines in Fig. 1).
- Unstratified Greenwood Z-test controls type I error.
- Z-tests with Borkowf's variance can be too conservative.

### Power of test:

- Data from COX model: The Cox regression performs best as expected, the Z-test with MH-weights and Borkowf is a reliable alternative.
- Data from ASD model: Z-test with IV-weights and Borkowf performs very well. Cox regression, unstratified Z-test with Greenwood and Z-test with MH-weights and Borkowf have almost identical power.



**Figure 1:** Power vs effect size for underlying COX model (left) and ASD model (right).

## CONCLUSIONS

- Z-tests for difference in survival at time  $t_0$**  are a **valuable alternative to Cox regression**, especially if the proportionality assumption does not hold. However, for small violations the Cox regression is still the model of choice.
- Greenwood variance can easily become zero** in small or extreme strata (Kaplan-Meier = 0 or 1)
- Z-tests with **Borkowf's variance control type I error** – stratified Z-test with Greenwood does not
- Mantel-Haenszel type weightings** seem promising (still assign a weight to strata with 0 variance)

## REFERENCES

- [1] Borkowf, C.B., 2005. A simple hybrid variance estimator for the Kaplan-Meier survival function. Statist. Med., 24, pp. 827-851.
- [2] Lachin, J. M., Biostatistical Methods. The Assessment of Relative Risks, 2nd ed. Wiley-Blackwell, 2011.
- [3] Greenland, S. and Robins, J.M., 1985. Estimation of a common effect parameter from sparse follow-up data. Biometrics, 41, pp. 55-68.

more information, including more scenarios:

