

# Statistical assessment of consistency of treatment effect in multiregional clinical trials

Kevin J Carroll



# Content and flow

- Background and key questions
- Design MRCT
- Analysis MRCT
- Examples
- Summary



# Background

- Modern drug development is impossible without the use of MRCTs
- However, there are challenges as reflected in ICH E5 relating to possible inconsistency of effects from region to region
  - 'intrinsic' = race, biology, genetics
  - 'extrinsic' = cultural, social, economic, medical practice, quality of care, experience in clinical trials, quality of trial conduct and monitoring
- With more recent regulatory guidelines from PMDA and EMA together with notable test cases, MRCTs subject to a lot of attention within and outside of statistics

# Content and flow

- Background and key questions
- **Design MRCT**
- Analysis MRCT
- Examples
- Summary



# Construct

- Trial designed to detect a true overall effect  $\theta$  with  $\alpha$  1-sided Type I error  $\alpha$  and power  $1-\beta$ . The information content is therefore

$$V = \frac{(z_\alpha + z_\beta)^2}{\theta^2}$$

- $i=1$  to  $r$  regions each with a fraction  $f_i$  patients. The estimated treatment effect estimate in  $i^{\text{th}}$  region is therefore

$$\hat{\theta}_i \sim N\left(\theta, \frac{1}{f_i V}\right) \text{ so that the overall estimate } \hat{\theta} \sim \sum_{i=1}^r f_i \hat{\theta}_i \sim N\left(\theta, \frac{1}{V}\right)$$

- However, if regional effects are considered as random then  $\theta_i \sim N(\theta, \tau^2)$

$$\text{so that } \hat{\theta}_i \sim N\left(\theta, \tau^2 + \frac{1}{f_i V}\right) \text{ and } \hat{\theta} \sim \sum_{i=1}^r f_i \hat{\theta}_i \sim N\left(\theta, \tau^2 \sum_{i=1}^r f_i^2 + \frac{1}{V}\right)$$

# Hung (2010)

- If, in truth,  $\tau^2 > 0$ , then Type I error will be inflated

$$1 - \Phi \left( z_\alpha \left( 1 + V\tau^2 \sum_{i=1}^r f_i^2 \right)^{-0.5} \right)$$

- V should be increased to

$$\tilde{V} = \left( \frac{\theta^2}{(z_\alpha + z_\beta)^2} - \tau^2 \sum_{i=1}^r f_i^2 \right)^{-1} \quad \text{so that} \quad \frac{V}{\tilde{V}} = 1 - \frac{\tau^2}{\theta^2} (z_\alpha + z_\beta)^2 \sum_{i=1}^r f_i^2$$

$$\text{and power} = \Phi \left( -z_\alpha + (z_\alpha + z_\beta) \left( 1 + V\tau^2 \sum_{i=1}^r f_i^2 \right)^{-0.5} \right)$$

- For  $V = \frac{N}{\sigma^2}$  sample size inflation is:

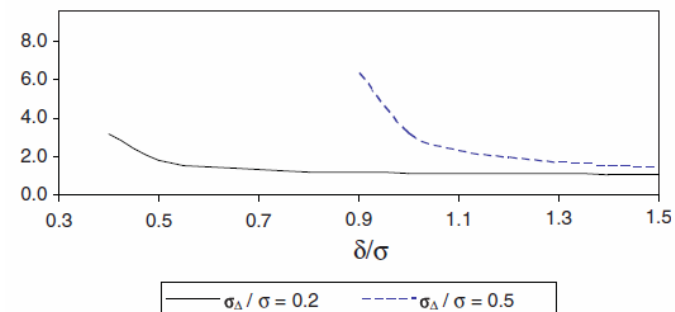
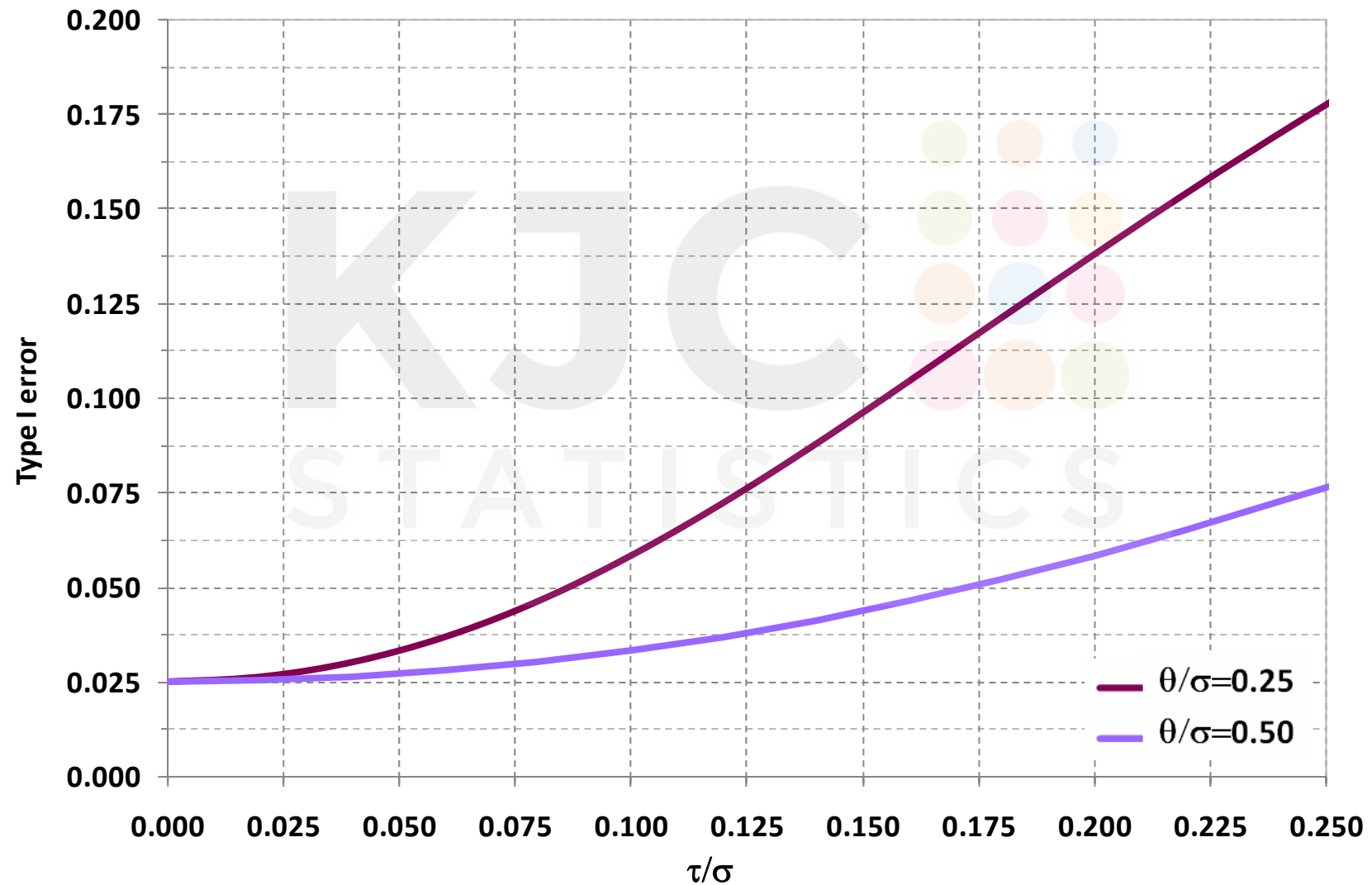


Figure 1. Sample size ratio  $N/N_0$  versus  $(\delta/\sigma)$ .

# Type I error inflation for $V=N/\sigma^2$ , $r=3$ , $f_i=1/3$ , $\beta=0.1$



# Issues and Questions (1)

- Ignoring  $\tau^2$  increases Type I error and reduces power.
- Should N in an MRCT be routinely increased, allowing for  $\tau^2$ ?
- In trial planning, what value for  $\tau^2$  (or  $\tau / \sigma$ ) should be assumed?  
On what is the choice based? Is it logical to assume  $\tau^2 > 0$ ?
- At what value of  $\tau / \sigma$  does a MRCT become
  - infeasible in terms of trial size
  - impossible to interpret clinically due to excessive between region variation in treatment effect?



## Issues and Questions (2)

- How do we define 'region'? Should there be a regulatory standard agreed to cover all trials ?
- What should be the allocation of N across regions / countries? How do we determine This?
- What is meant by 'consistency' ? How do we define this? How do we assess it? What is the value and role of routine homogeneity testing of regional results? And graphical methods?
- Should a random effects analysis be the standard in MRCTs? What are the consequences if so?

# What allocation of patient to regions?

- If we want to observe a positive treatment effect in each region, then

$$\Pr(\hat{\theta}_i > 0, \forall i = 1 \text{ to } r) = \prod_{i=1}^r \Phi\{\sqrt{f_i}(z_\alpha + z_\beta)\}$$

- But need also to observe an overall significant effect, so then we want

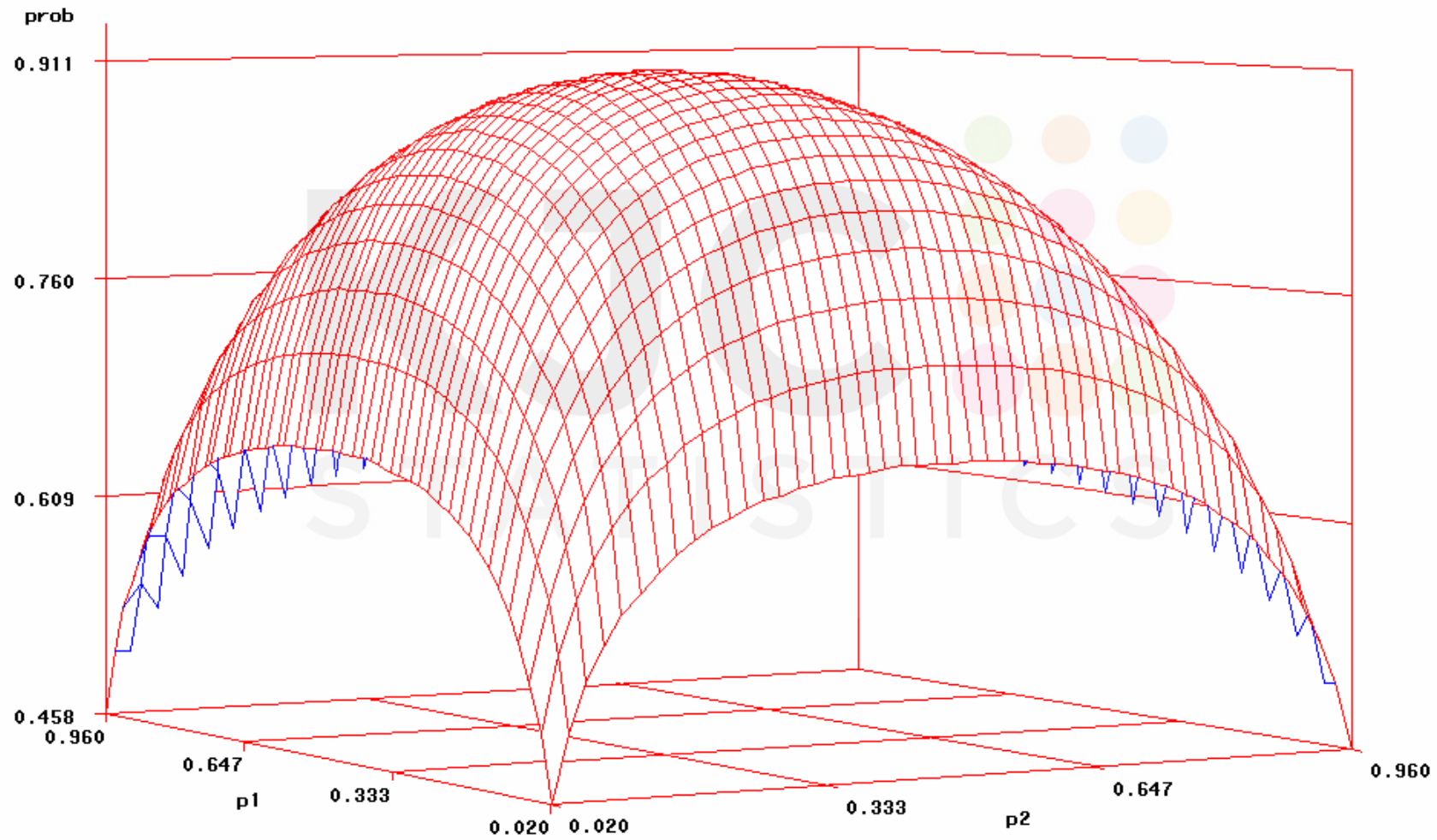
- $\Pr(\hat{\theta}_i > 0, \forall i = 1 \text{ to } r \cap \hat{\theta} > z_\alpha V^{-0.5})$  where  
 $(\hat{\theta}, \hat{\theta}_1, \dots, \hat{\theta}_r) \sim N(\underline{\theta}, \underline{\Lambda})$

$$\underline{\Lambda} = \begin{bmatrix} \frac{1}{v} & \frac{1}{v} & \cdot & \frac{1}{v} \\ \frac{1}{v} & \frac{1}{f_1 v} & 0 & 0 \\ \cdot & 0 & \cdot & 0 \\ \frac{1}{v} & 0 & \cdot & \frac{1}{v} \end{bmatrix}$$

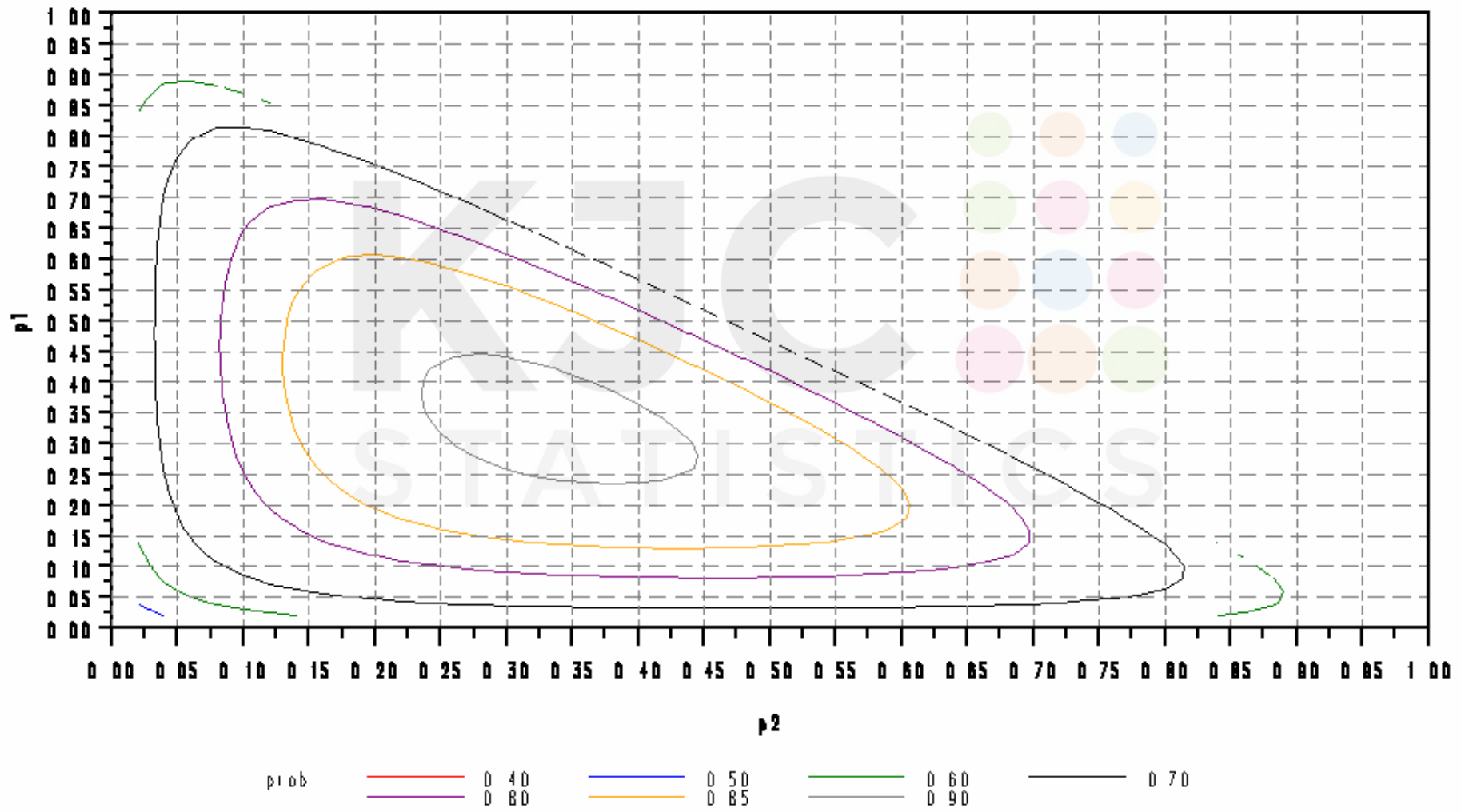
- Chuang (2008) considers the most typical MRCT with 3 regions, to reflect US, EU and Japan concluding the smallest region could be as low as 15% and yet still achieve = 0.80.

$$\Pr(\hat{\theta}_i > 0, \forall i = 1 \text{ to } r \cap \hat{\theta} > z_\alpha V^{-0.5})$$

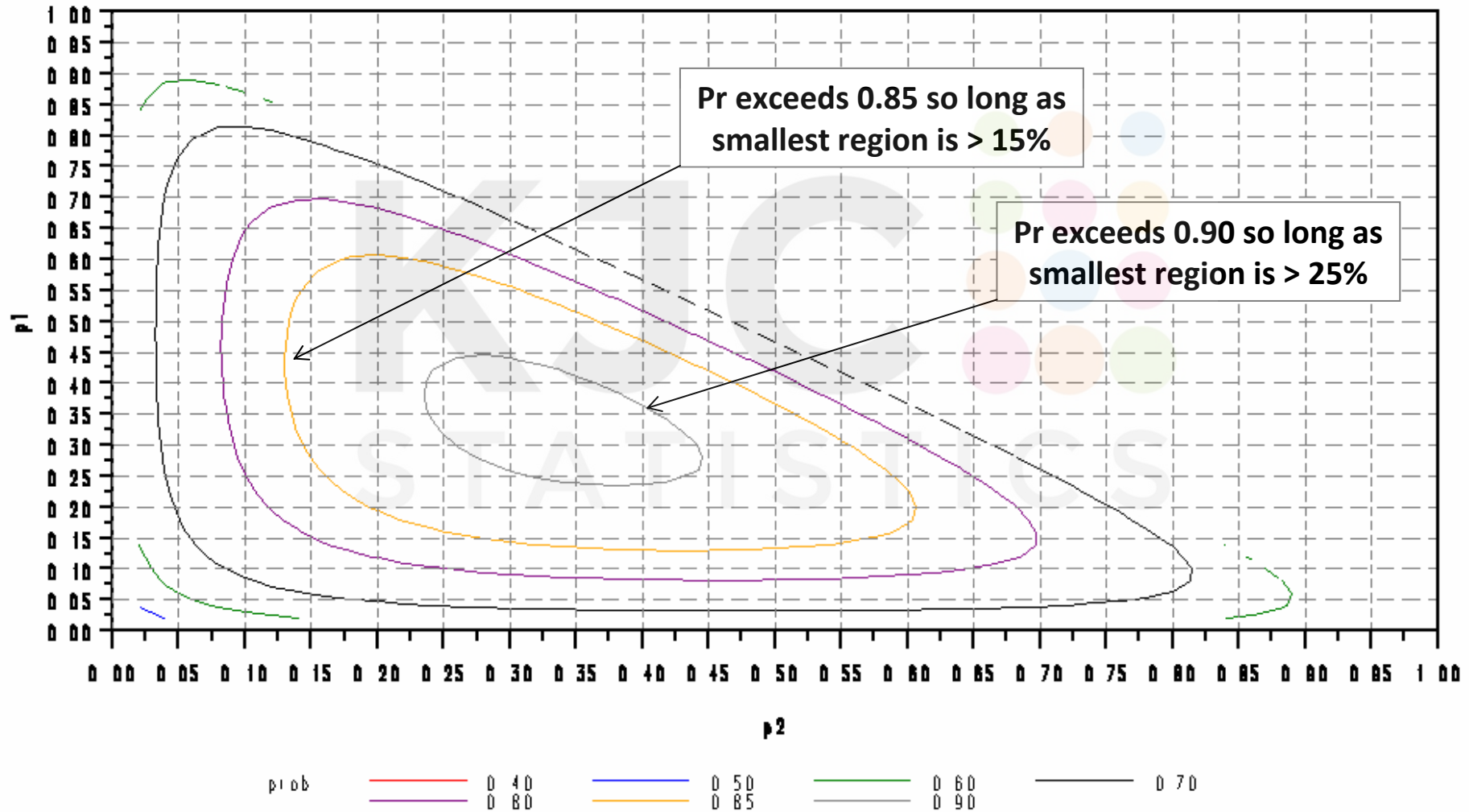
# Pr observing a positive treatment effect in each region



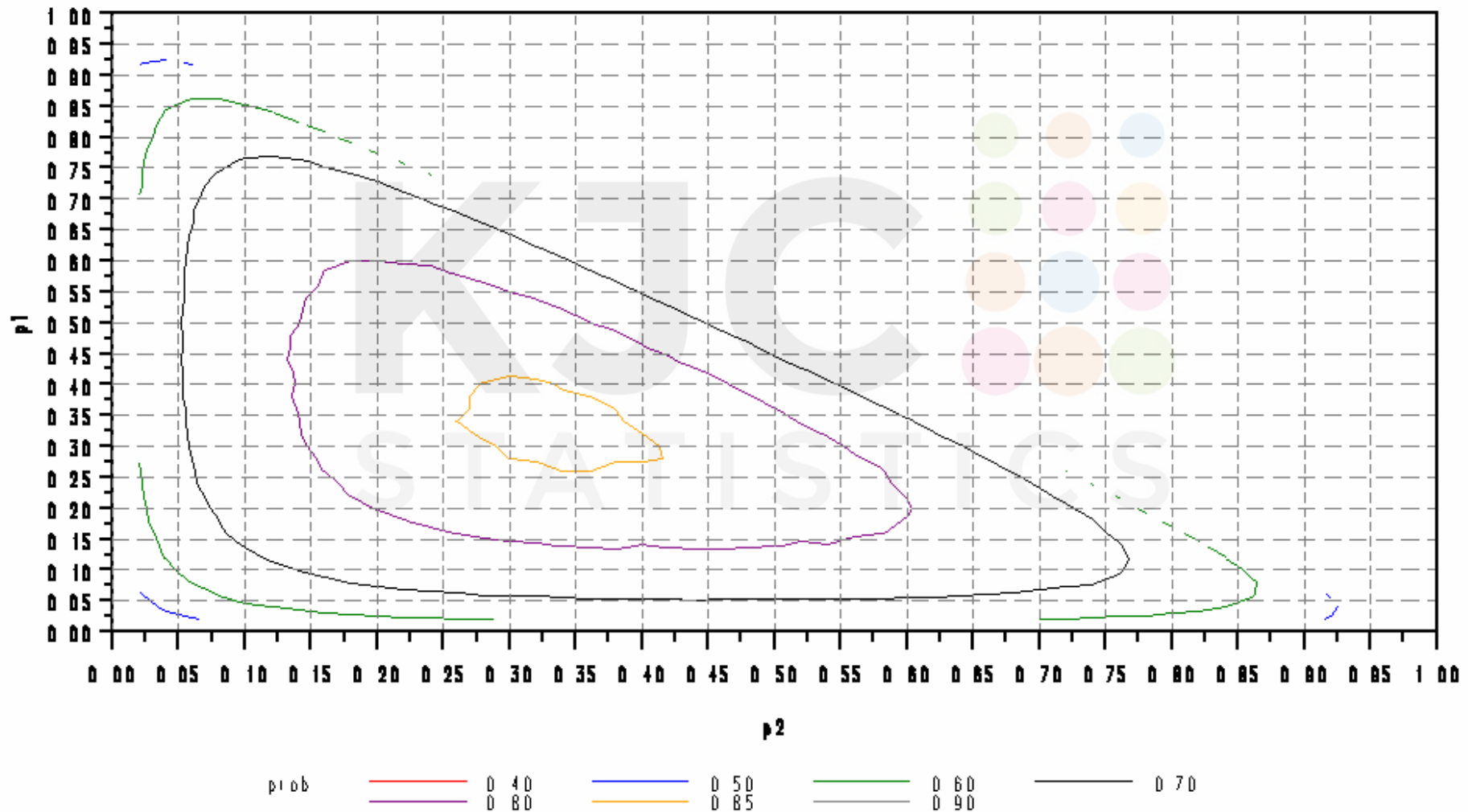
# Pr observing a positive treatment effect in each



# Pr observing a positive treatment effect in each

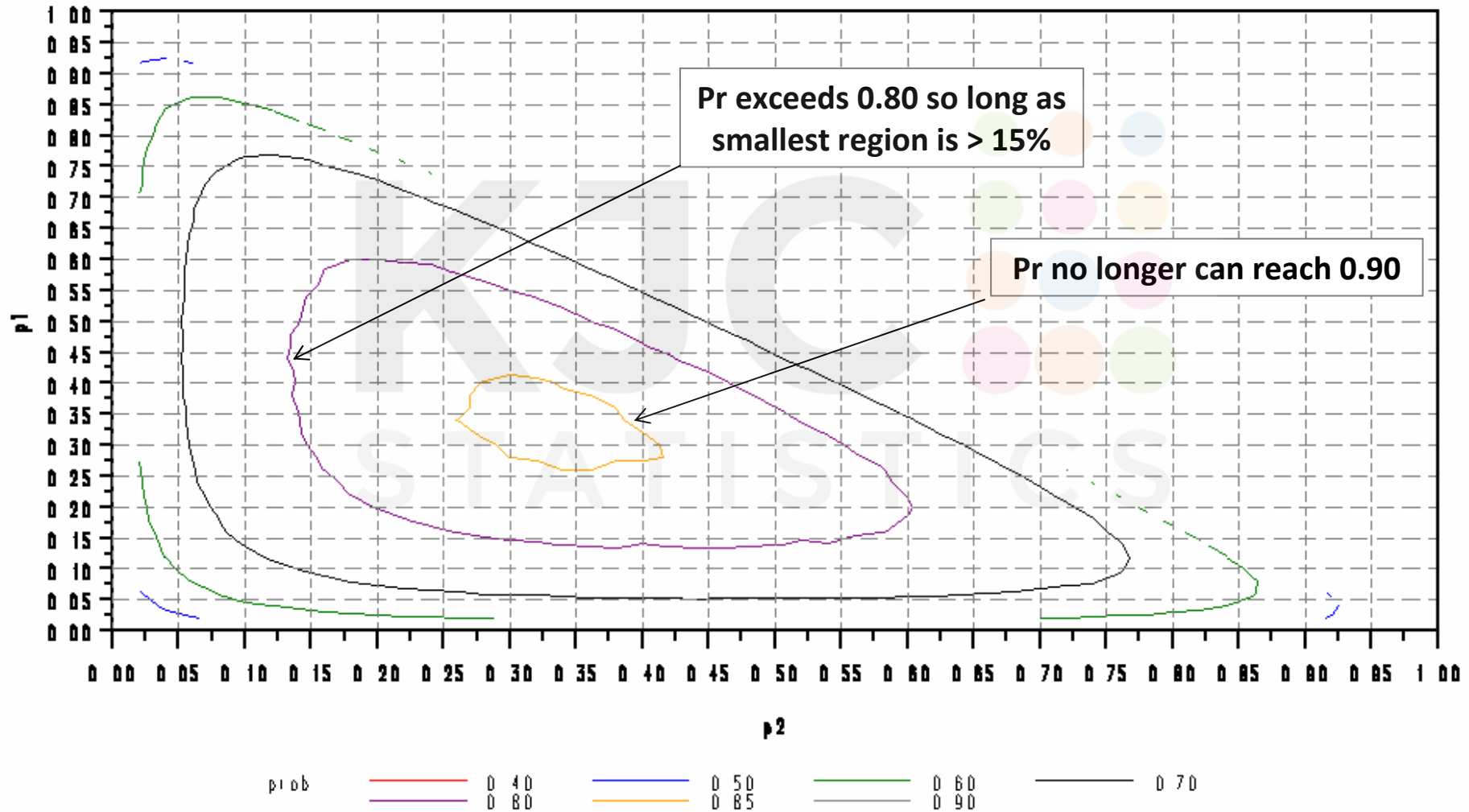


# Pr observing a positive treatment effect in each region and $p < 0.05$ for the overall effect



Based on 10,000 simulations for each combination of  $p_1$  and  $p_2$

# Pr observing a positive treatment effect in each region and $p < 0.05$ for the overall effect



Based on 10,000 simulations for each combination of  $p_1$  and  $p_2$

# Looking at the 'reference' region of interest<sup>1</sup>

- Denote  $\hat{\theta}_{\text{ref}} \sim N\left(\theta, \frac{1}{fV}\right)$  and  $\hat{\theta}_{\text{nref}} \sim N\left(\theta, \frac{1}{(1-f)V}\right)$  then consistency between the reference region and non-reference regions may be defined in several ways:

1.  $\hat{\theta}_{\text{ref}} - 0.5\hat{\theta} > 0$

2.  $\hat{\theta}_{\text{ref}} - 0.5\hat{\theta}_{\text{nref}} > 0$

3.  $\left(\hat{\theta}_{\text{ref}} - \hat{\theta}_{\text{nref}}\right) + \frac{\theta}{2} > 0$

4.  $\hat{\theta}_{\text{ref}} - \left(\hat{\theta} - z_{\alpha} V^{-0.5}\right) > 0$

5.  $\hat{\theta}_{\text{ref}} > 0$  and  $\hat{\theta}_{\text{nref}} > 0$

- Denoting criteria as  $C_j$  then  $\text{pr}(\text{meeting criteria and achieving } \hat{\theta} > z_{\alpha} V^{-0.5})$  can be determined by noting

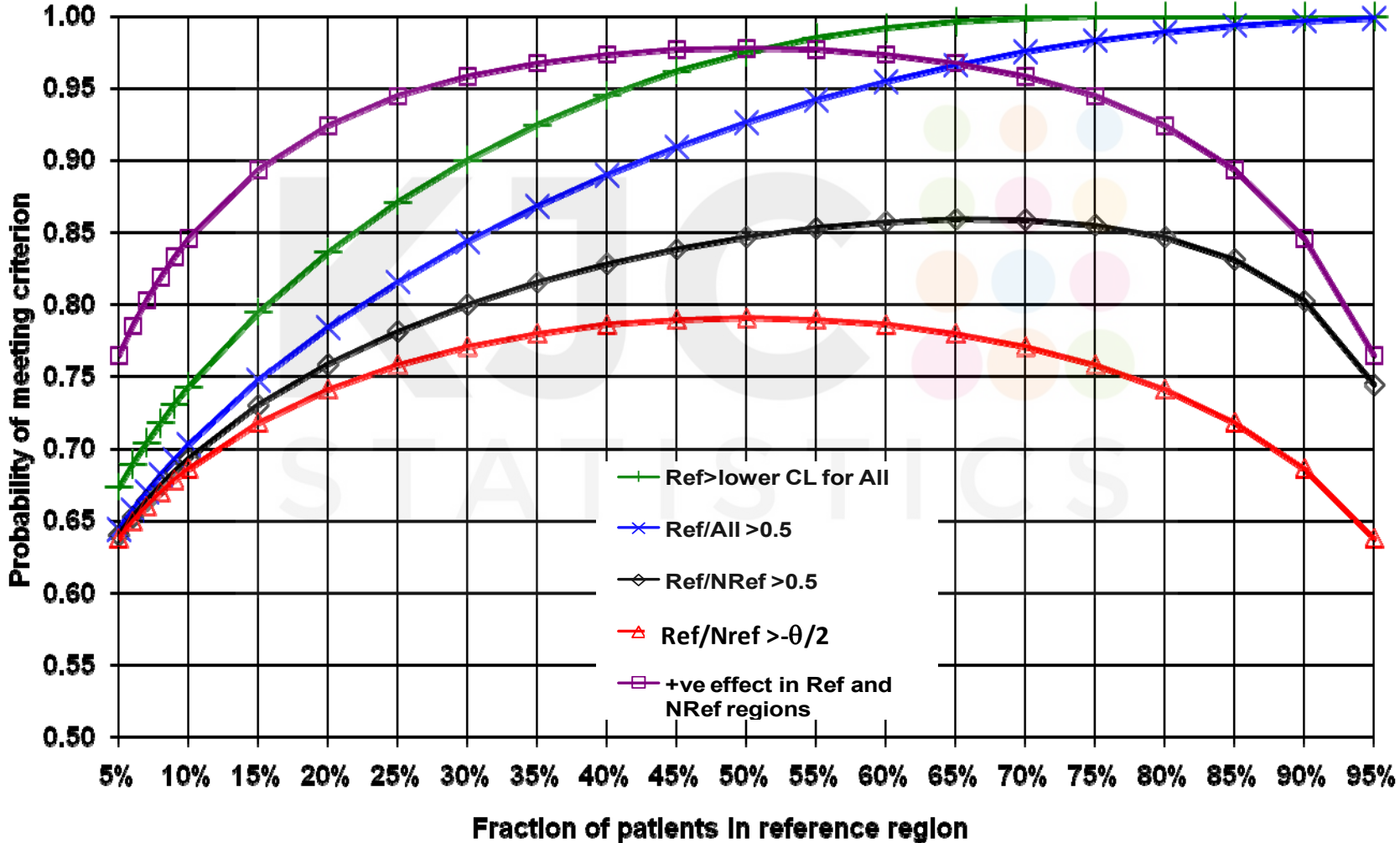
$$\begin{pmatrix} C_1 \\ \hat{\theta} \end{pmatrix} \sim N \begin{pmatrix} \theta(1-k), \frac{1}{fV} + \frac{k(k-2)}{V}, \frac{1-k}{V} \\ \theta, \frac{1-k}{V}, \frac{1}{V} \end{pmatrix}$$

$$\begin{pmatrix} C_2 \\ \hat{\theta} \end{pmatrix} \sim N \begin{pmatrix} \theta(1-k), \frac{1}{fV} + \frac{k^2}{V(1-f)}, \frac{1-k}{V} \\ \theta, \frac{1-k}{V}, \frac{1}{V} \end{pmatrix}$$

$$\begin{pmatrix} C_4 \\ \hat{\theta} \end{pmatrix} \sim N \begin{pmatrix} z_{\alpha} V^{-0.5}, \frac{(1-f)}{fV}, 0 \\ \theta, 0, \frac{1}{V} \end{pmatrix}$$



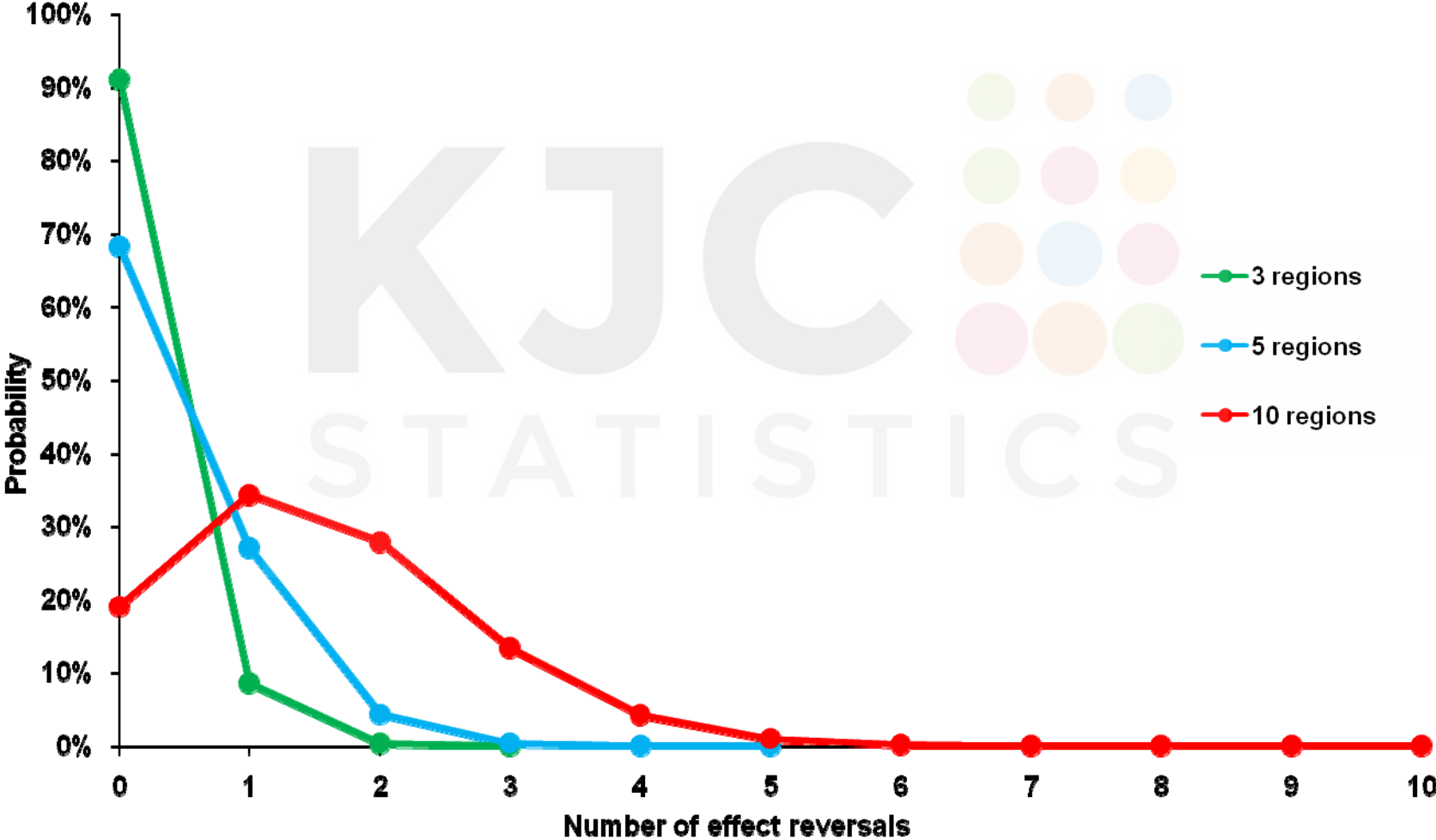
# Pr(meeting criteria 1 to 5)



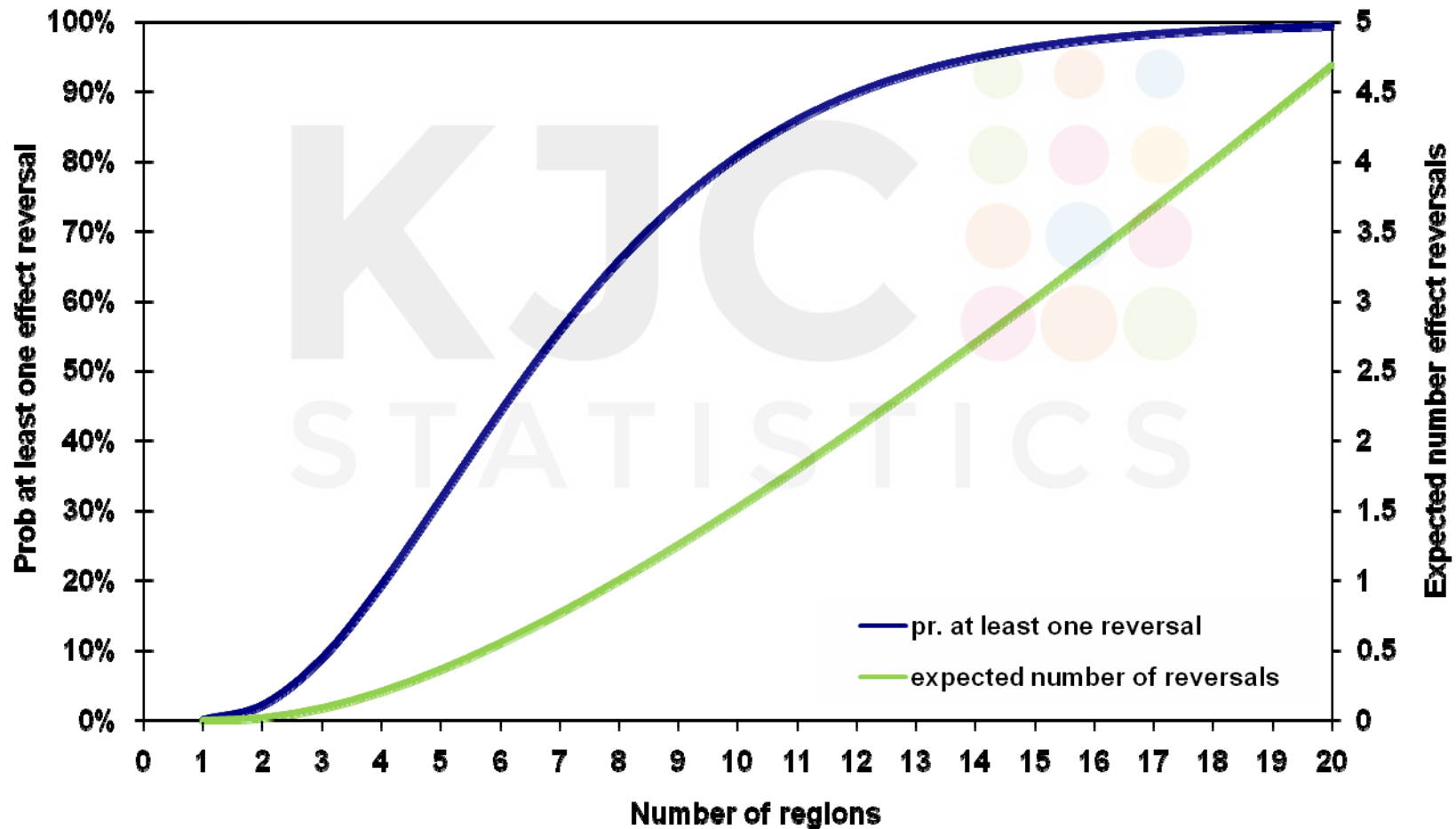
## So what?

- Of these possible alternative criteria, only criterion #4 represents some improvement over the MHLW guideline criterion in terms of the fraction of reference region patients required in a MRCT
- Criterion #4 requires around 1/3 fewer patients than Criterion #1 for about the same overall power.

“Effect Reversals”, where the treatment effect is positive overall but numerically negative in some regions are to be expected in a large multiregional trials



In a trial with 90% (80%) power and 7 (6) regions, the probability of at least one effect reversal >50%.



# Content and flow

- Background and key questions
- Design MRCT
- **Analysis MRCT**
- Examples
- Summary



# Random effects analysis – good idea?

- If regional effects are considered as random then  $\theta_i \sim N(\theta, \tau^2)$  and

$$\hat{\theta}_i \sim N\left(\theta, \tau^2 + \frac{1}{f_i V}\right) \quad \text{with} \quad Q = v \sum_{i=1}^r f_i (\hat{\theta}_i - \hat{\theta})^2 \sim \chi_{r-1}^2$$

- Then  $\hat{\tau}^2 = \frac{Q - (r-1)}{1 - \sum_{i=1}^r f_i^2}$  and  $w_i = \left(\hat{\tau}^2 + \frac{1}{f_i V}\right)^{-1}$  so that

$$\hat{\theta}_{RE} \sim \sum_{i=1}^r w_i \hat{\theta}_i \sim N\left(\theta, \left(\sum_{i=1}^r w_i\right)^{-1}\right)$$

- Example, Chen (2010):

$$\hat{\theta} = 0.89 \text{ 95\% CI } (0.79, 0.99), p=0.037$$

$$\hat{\theta}_{RE} = 0.91 (0.76, 1.08), p=0.29$$

Table I. PURSUIT: Efficacy by Region.

	N	Odds ratio	95% confidence interval
Overall	10948	0.89	(0.79, 0.99)
North America	4358	0.75	(0.63, 0.91)
Western Europe	4243	0.92	(0.77, 1.11)
Latin America	585	1.03	(0.60, 1.76)
Eastern Europe	1762	1.09	(0.85, 1.39)

N: sample size. Odds ratio: eptifibatide vs placebo, lower the better [19].

# Random effects analysis – good idea?

- Example, Chen (2010):

$$\hat{\tau}^2 = 0.016, \hat{\tau}/\hat{\sigma} = 1.06$$

$$\chi^2 = 6.3, p=0.096, I^2=0.52$$

$$\hat{\theta} = 0.89 \text{ CI } (0.79, 0.99), p=0.037$$

$$\hat{\theta}_{RE} = 0.91 \text{ CI } (0.76, 1.08), p=0.29$$

Table I. PURSUIT: Efficacy by Region.

	N	Odds ratio	95% confidence interval
Overall	10948	0.89	(0.79, 0.99)
North America	4358	0.75	(0.63, 0.91)
Western Europe	4243	0.92	(0.77, 1.11)
Latin America	585	1.03	(0.60, 1.76)
Eastern Europe	1762	1.09	(0.85, 1.39)

N: sample size. Odds ratio: eptifibatide vs placebo, lower the better [19].

- $\hat{\theta}_{RE} = 0.91$  (0.68, 1.21),  $p=0.36$  if Follmann<sup>1</sup> adjustment applied which means using a t value on k-1 df for  $\hat{\theta}_{RE}$  as opposed to z value. This has been advocated in the recent past by FDA for meta-analyses .
- If applied this has some odd consequences for MRCTs.

	fi	RR	Lower 95% CL	Upper 95% CL	p-value
Region 1	0.33	0.80	0.65	0.99	0.0416
Region 2	0.33	0.80	0.65	0.99	0.0416
Region 3	0.33	0.80	0.65	0.99	0.0416
All		0.80	0.71	0.91	0.0004
RE		0.80	0.61	1.05	0.0718

# Content and flow

- Background and key questions
- Design MRCT
- Analysis MRCT
- **Examples**
- Summary



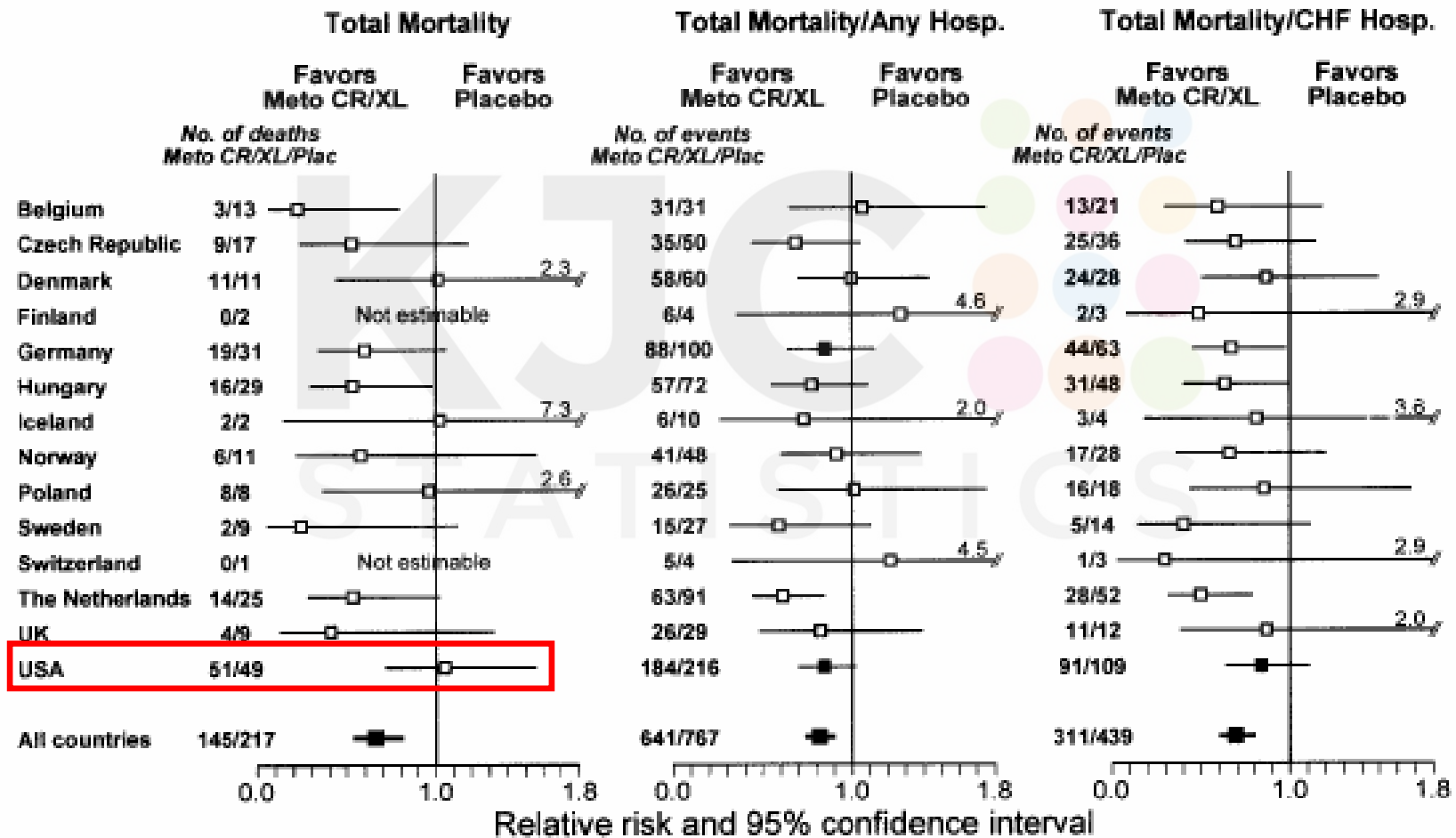


## MERIT-HF

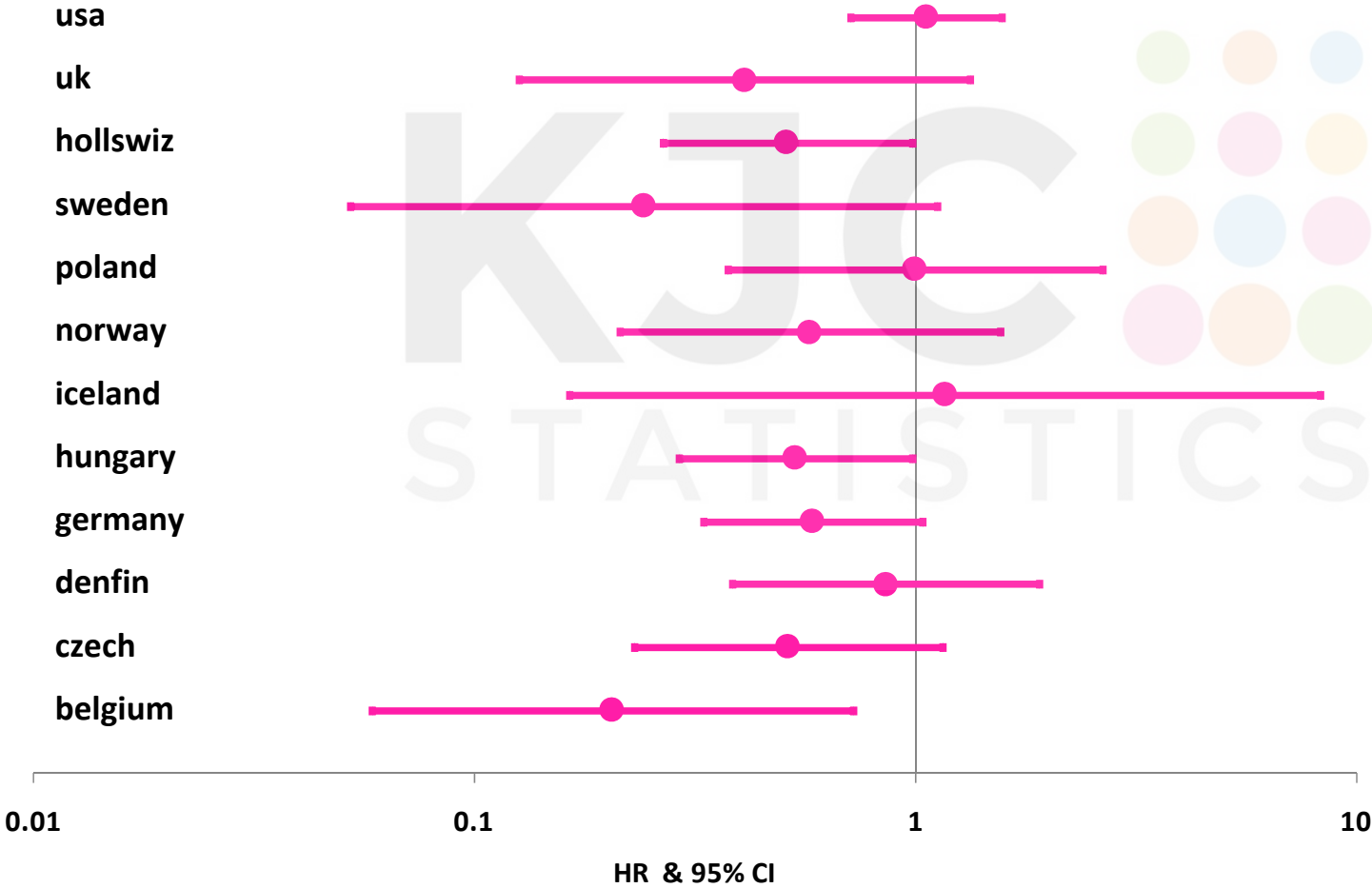
- Examined metoprolol CR/XL vs placebo in patients chronic heart failure and decreased ejection fraction
- Randomised 3991 across 14 countries
- The study stopped early on the recommendation of the IDMC. All-cause mortality was lower in the metoprolol CR/XL group than in the placebo group (145 [7.2%] vs 217 deaths [11.0 %]) RR = 0.66 [95% CI 0.53–0.81];  $p=0.00009$ .
- However....

# MERIT-HF: Inconsistent regional effects?

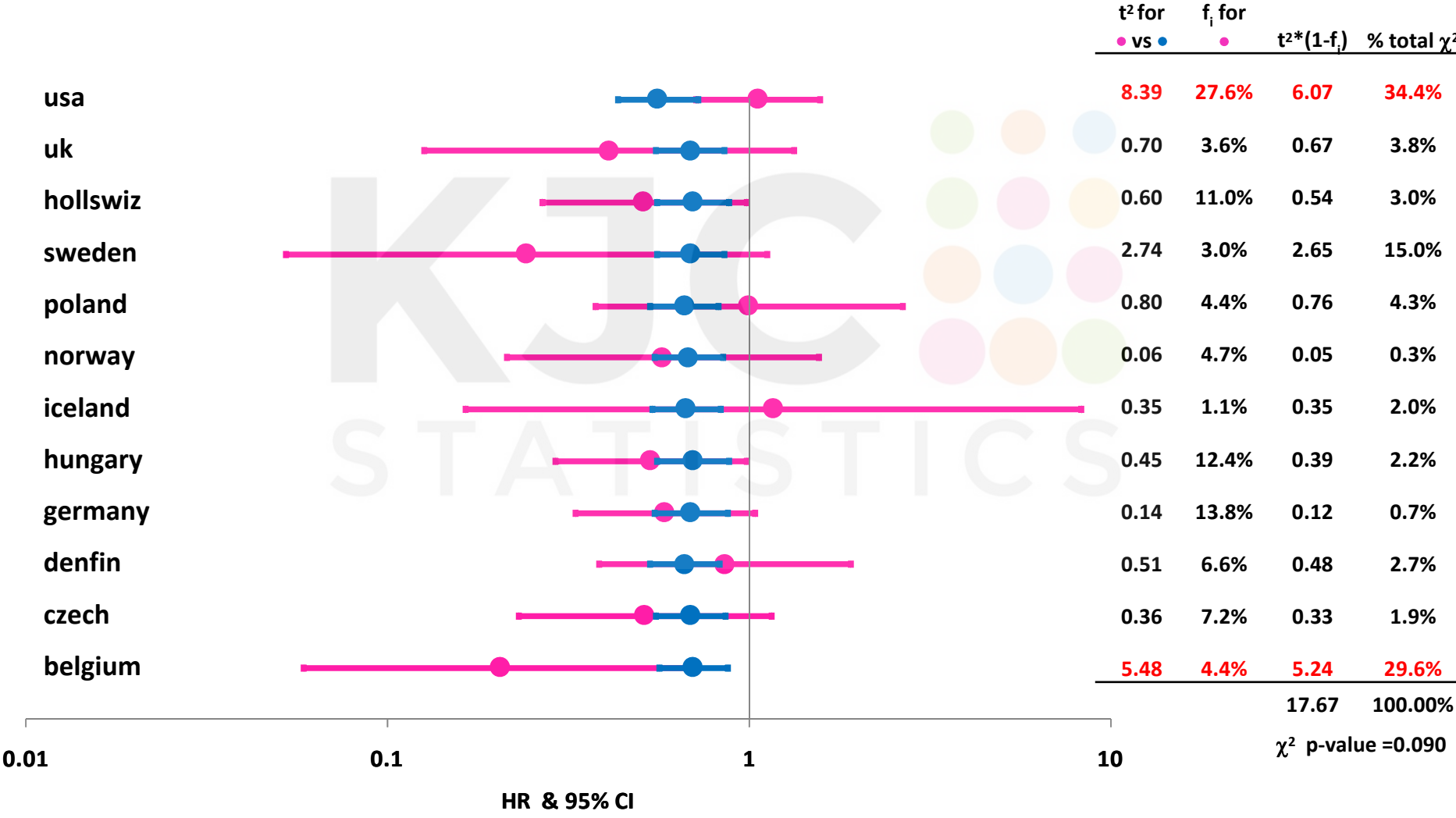
## All Patients Randomized



# MERIT-HF: Mortality by region

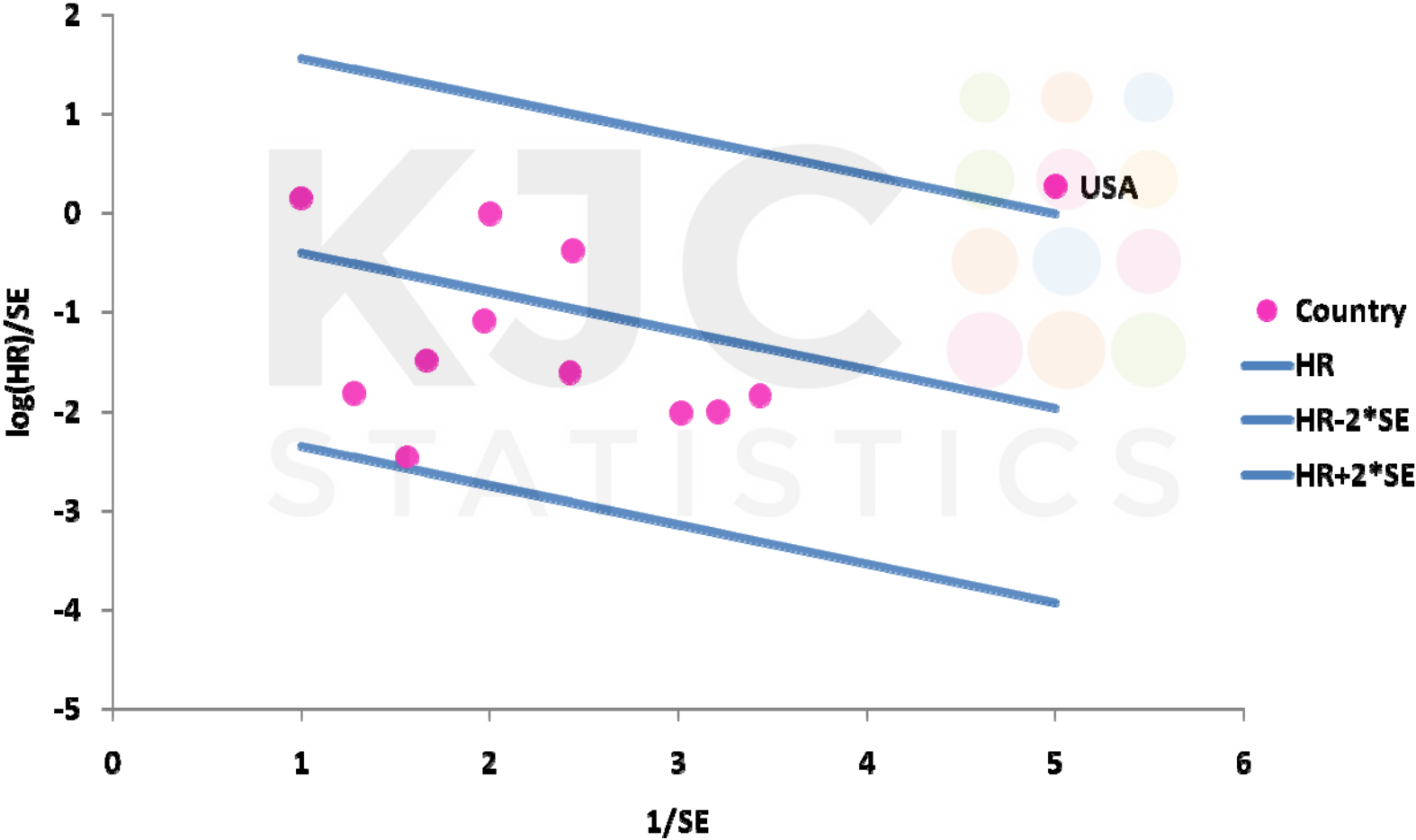


# MERIT-HF: Mortality by region vs all other regions



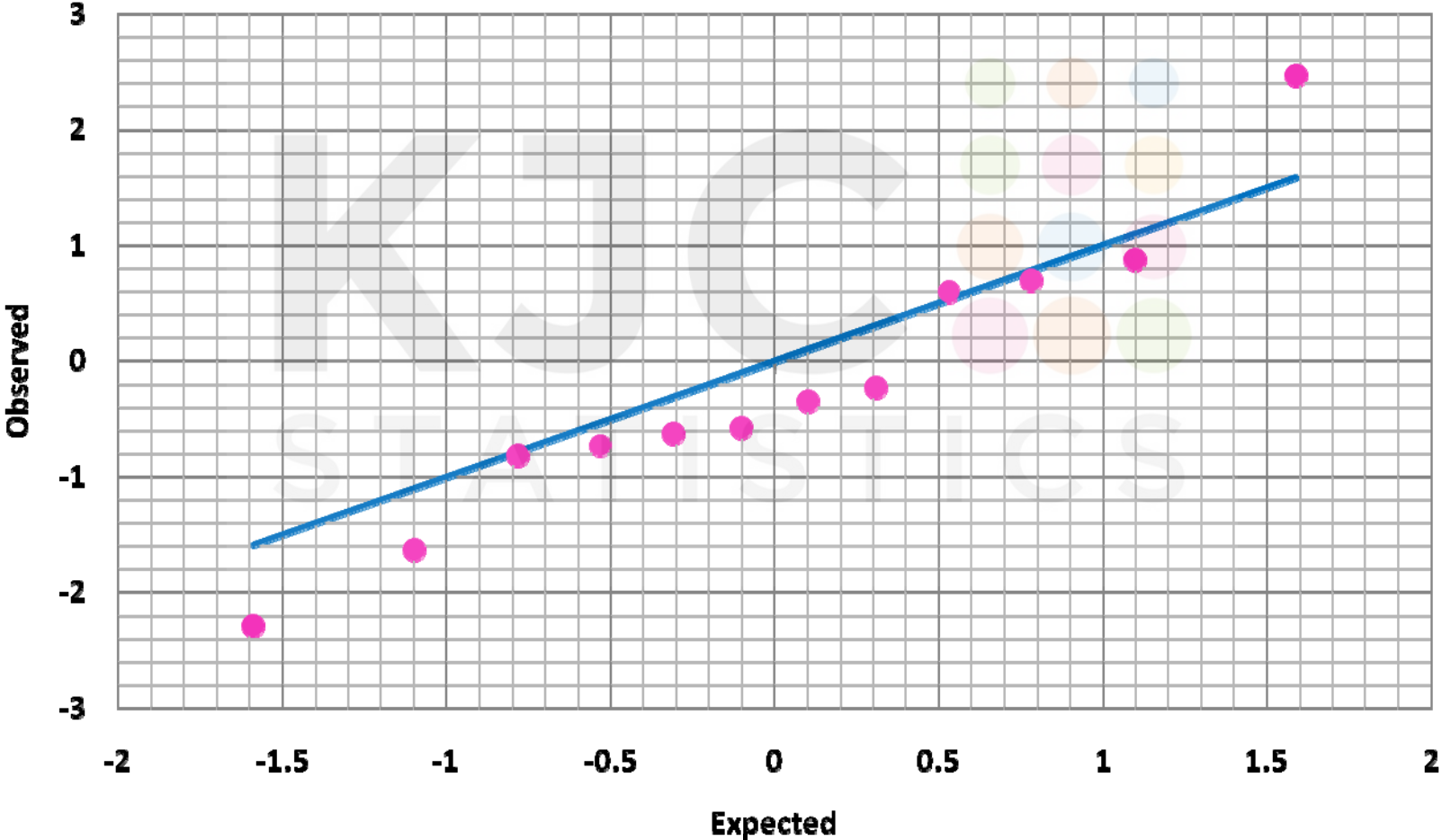
# MERIT-HF: Mortality:

HR=0.65, 95% CI (0.49. 0.86),  $\tau^2=0.038$ ,  $s^2=0.14$ ,  $I^2=0.21$



# MERIT-HF: Mortality:

HR=0.65, 95% CI (0.49, 0.86),  $\tau^2=0.038$ ,  $s^2=0.14$ ,  $I^2=0.21$



# MERIT-HF: Mortality vs Primary

## Mortality

Country	Drug		placebo		vi	xi
	Events	N	Events	N		
belgium	3	68	13	66	0.3804	-1.4962
czech	9	123	17	124	0.1537	-0.6279
Denm/Fin	11	161	13	164	0.1555	-0.1486
germany	19	252	31	247	0.0769	-0.5096
hungary	16	211	29	212	0.0875	-0.5900
iceland	2	19	2	22	0.9019	0.1466
norway	6	97	11	105	0.2377	-0.5269
poland	8	102	8	102	0.2304	0.0000
sweden	2	39	9	46	0.5637	-1.3390
holl/swiz	14	299	26	291	0.1031	-0.6462
uk	4	87	9	83	0.3376	-0.8580
usa	51	532	49	539	0.0363	0.0531

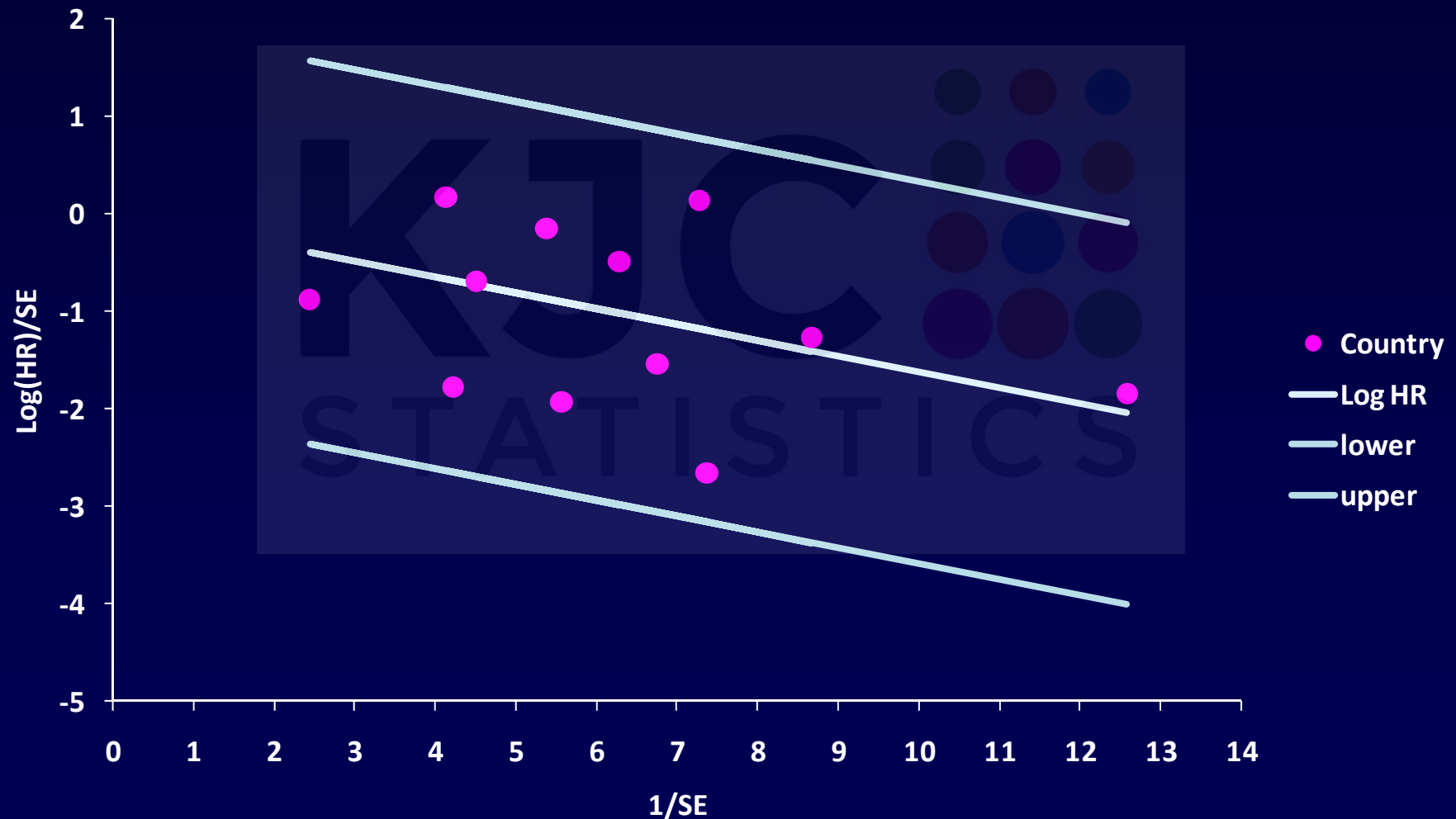
Given the overall result and based on distribution of events across countries, expect 2 countries to show effects > 0

## Primary (mortality + hospitalisation)

Country	Drug		placebo		vi	xi
	Events	N	Events	N		
belgium	31	68	31	66	0.0347	-0.0299
czech	35	123	50	124	0.0324	-0.3486
Denm/Fin	64	161	64	164	0.0189	0.0185
germany	88	252	100	247	0.0134	-0.1479
hungary	57	211	72	212	0.0220	-0.2289
iceland	6	19	10	22	0.1686	-0.3642
norway	41	97	48	105	0.0254	-0.0784
poland	26	102	25	102	0.0589	0.0392
sweden	15	39	27	46	0.0563	-0.4227
holl/swiz	68	299	95	291	0.0185	-0.3615
uk	26	87	29	83	0.0494	-0.1563
usa	184	532	216	539	0.0063	-0.1473

Given the overall result and based on distribution of events across countries, expect 2 countries to show effects > 0

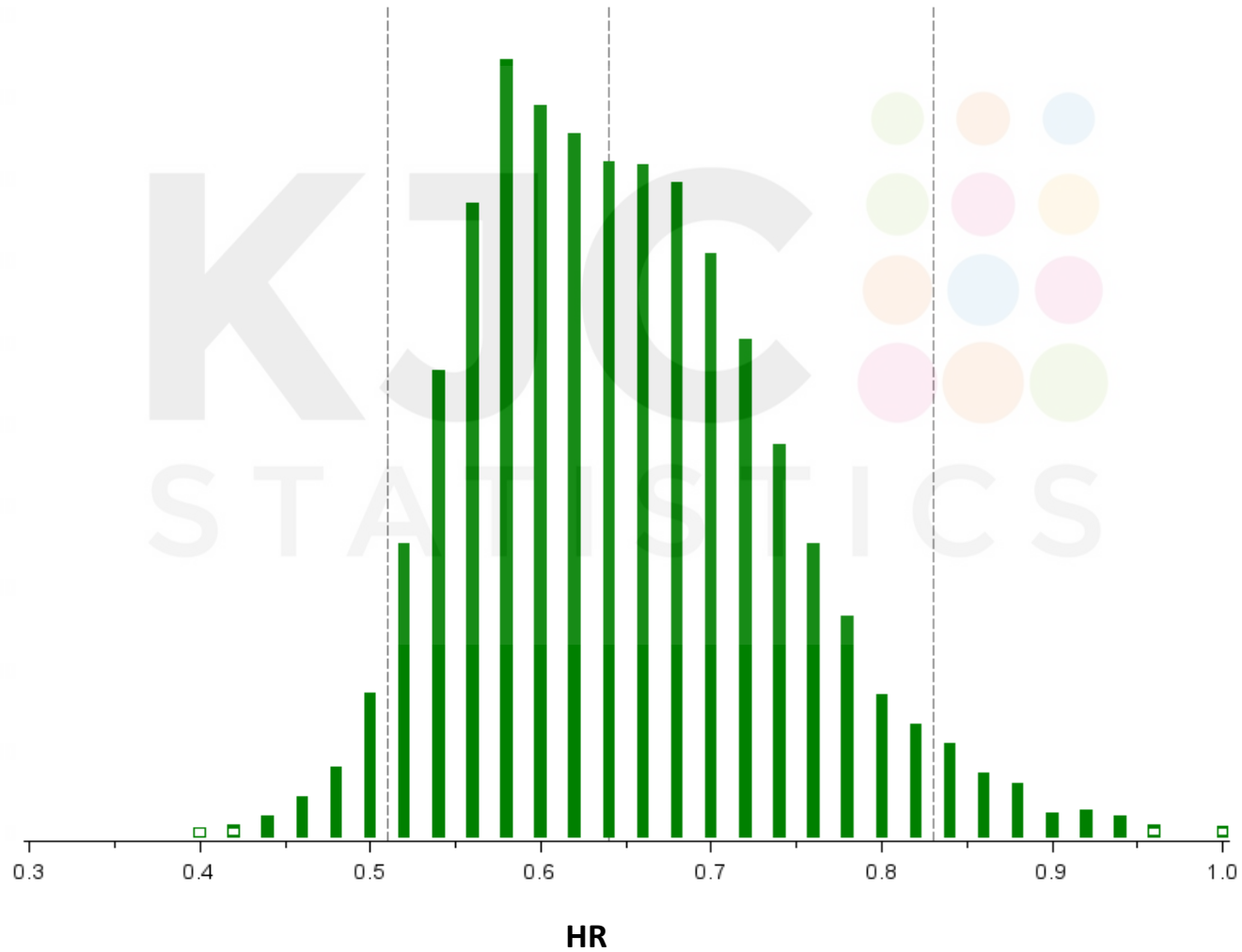
MERIT-HF: Primary (mortality + hospitalisaion):  
HR=0.85 95% CI (0.77. 0.93),  $\tau^2=0$ ,  $\sigma^2=0.024$ ,  $I^2=0$





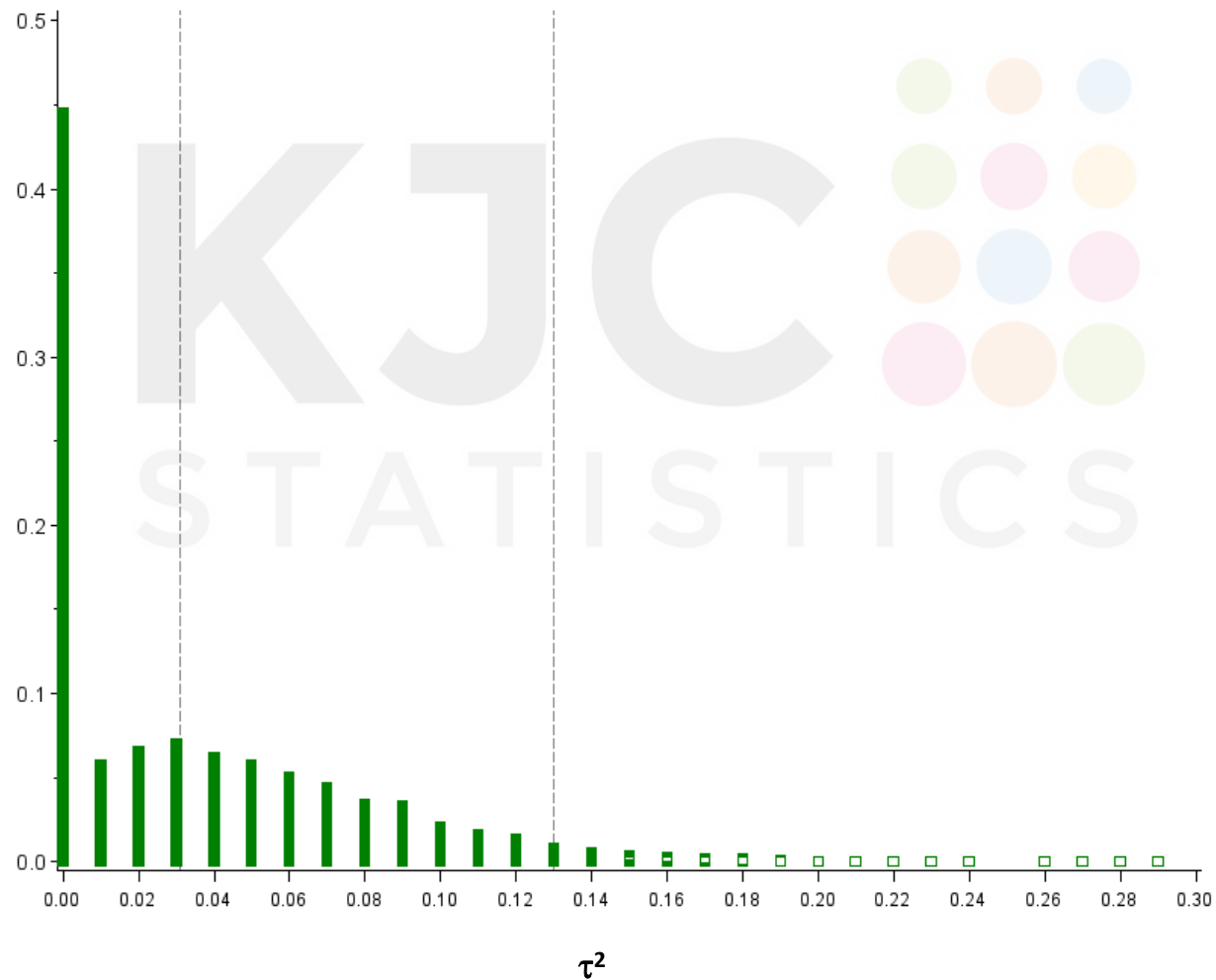
# MERIT-HF:

Mortality: HR=0.64 95% CI (0.51. 0.83) by bootstrap

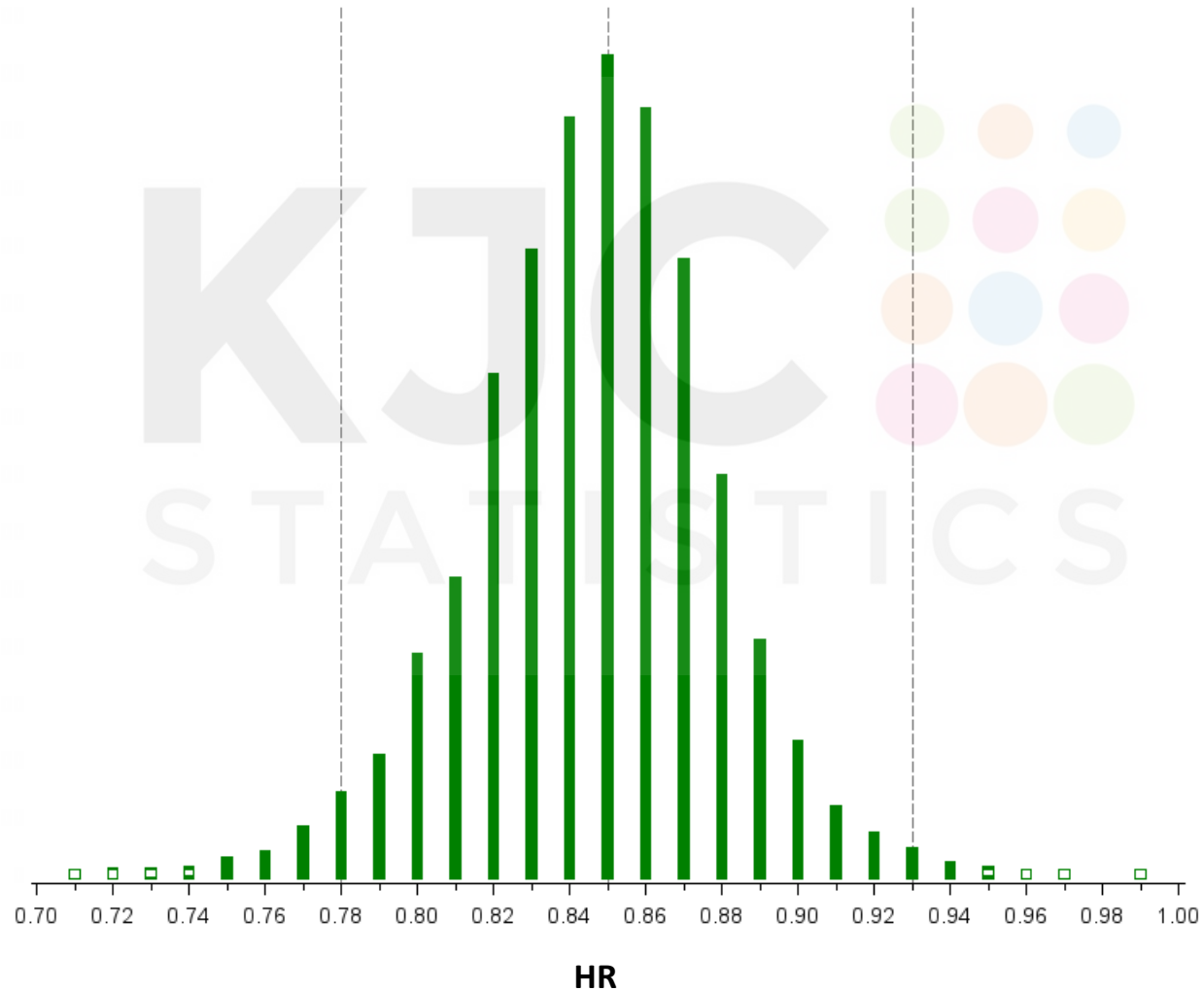


# MERIT-HF:

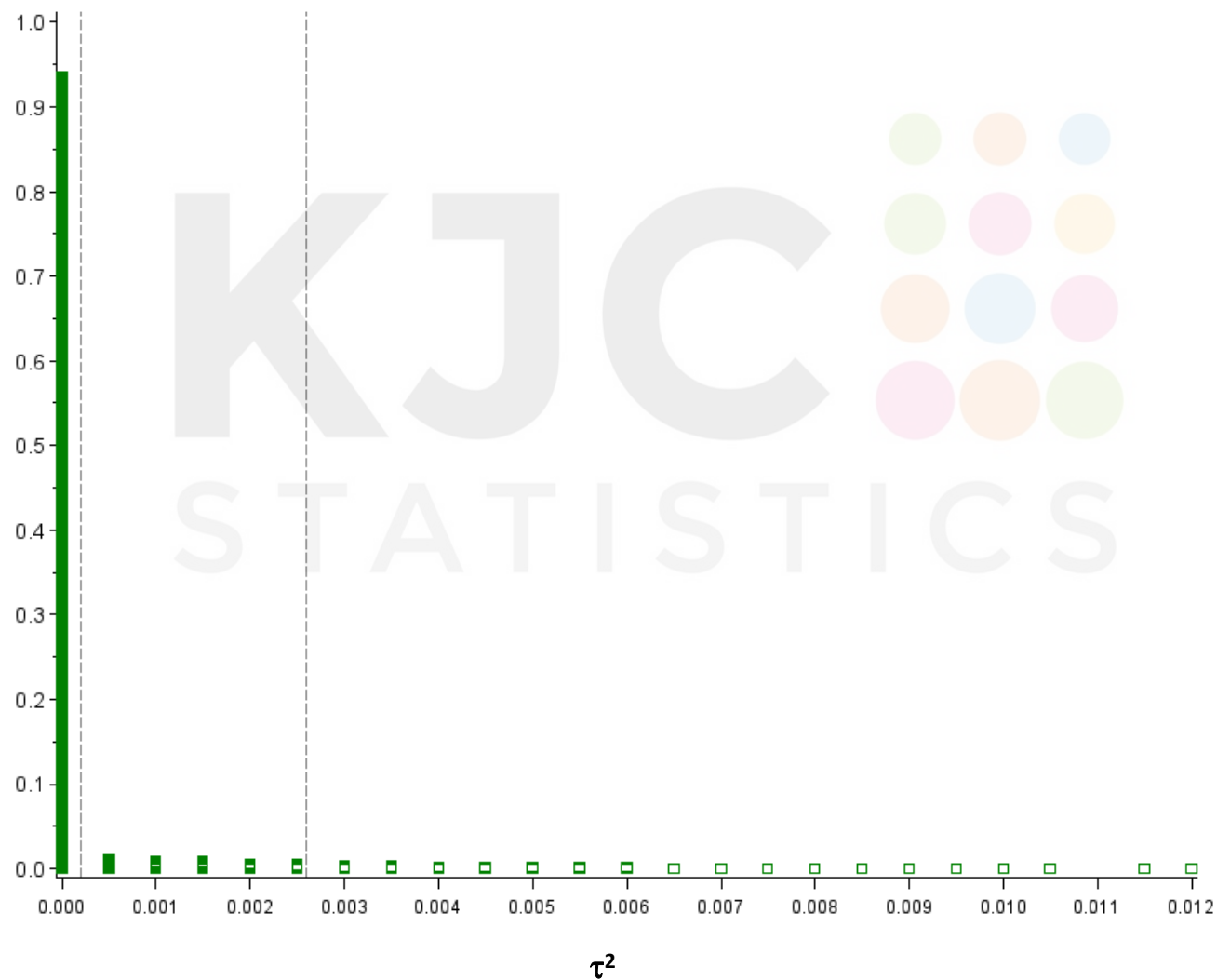
Mortality:  $\tau^2=0.031$  95% CI (0. 0.13)



**MERIT-HF: Primary (mortality + hospitalisation):  
HR=0.85 95% CI (0.79 to 0.90) by bootstrap**



# MERIT-HF: Primary (mortality + hospitalisation): $\tau^2=0$ , 95% CI (0, 0.0025)



## Wedel and DeMets, (Am Heart J 2001;142:502-11.)

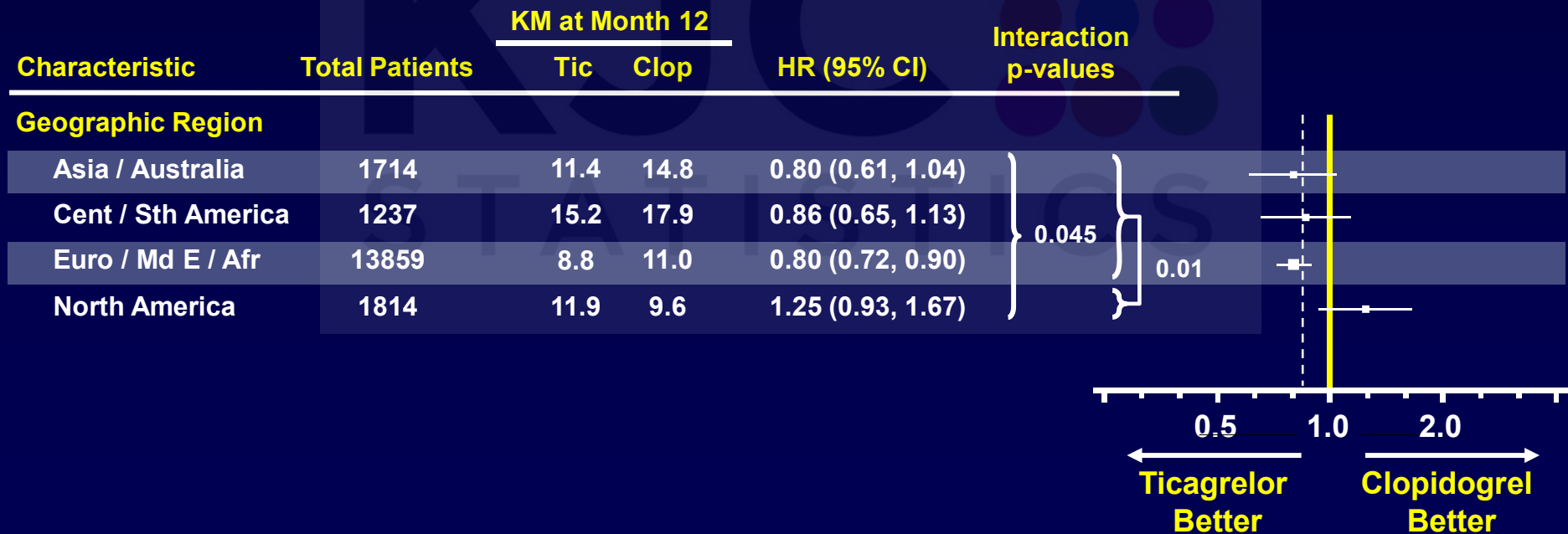
- Just as we must be extremely cautious in over-interpreting positive effects in subgroups, even those that are predefined, we must also be cautious in focusing on subgroups with an apparent neutral or negative trend.
- We should examine subgroups to obtain a general sense of consistency, which is clearly the case in MERIT-HF.
- We should expect some variation of the treatment effect around the overall estimate as we examine a large number of subgroups because of small sample size in subgroups and chance.
- Thus the best estimate of the treatment effect on total mortality for any subgroup is the estimate of the hazard ratio for the overall trial.

# PLATO

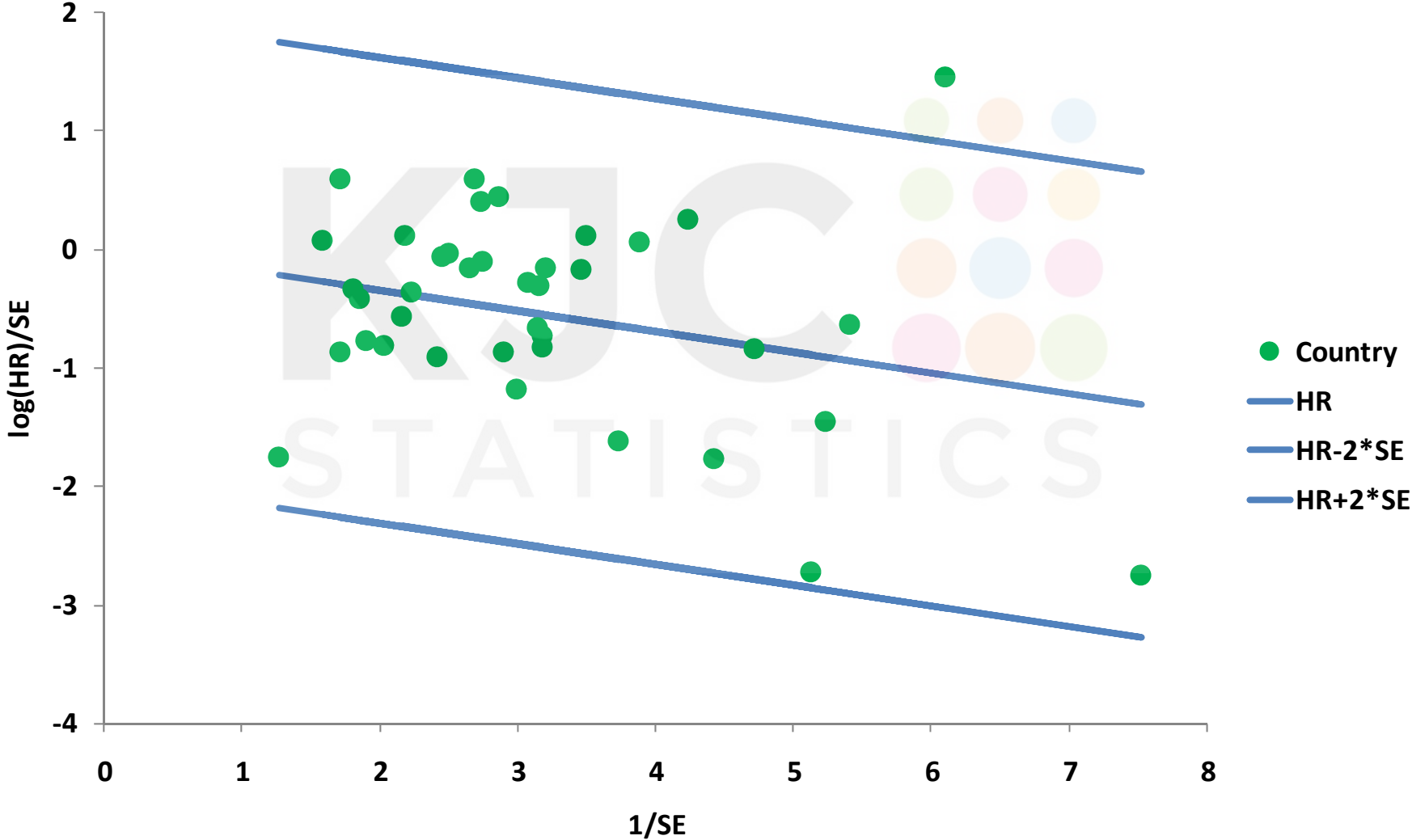
- Randomized double-blind study comparing BRILINTA (N=9333) to clopidogrel (N=9291), both given in combination with aspirin, in patients with acute coronary syndromes.
- Primary endpoint was time to first occurrence of CV death, MI or stroke.
- Randomisation across 41 countries.
- Primary endpoint met for BRILINTA 9.8% vs 11.7% events HR = 0.84 95% CI 0.77–0.92]; p=0.0003.
- Benefit also seen in overall mortality 4.5% vs 5.9% events HR = 0.78 95% CI 0.69–0.89]; p=0.0003.

# PLATO: Ticagrelor Effect Apparently Inconsistent Across Geographic Regions

- 31 pre-specified subgroup tests conducted for consistency
- No  $\alpha$ -level adjustment for multiplicity
- Indication of qualitatively different outcomes by region
- Results in NA appear to be driven by US: HR 1.27 (0.92, 1.75)



# Galbraith plot





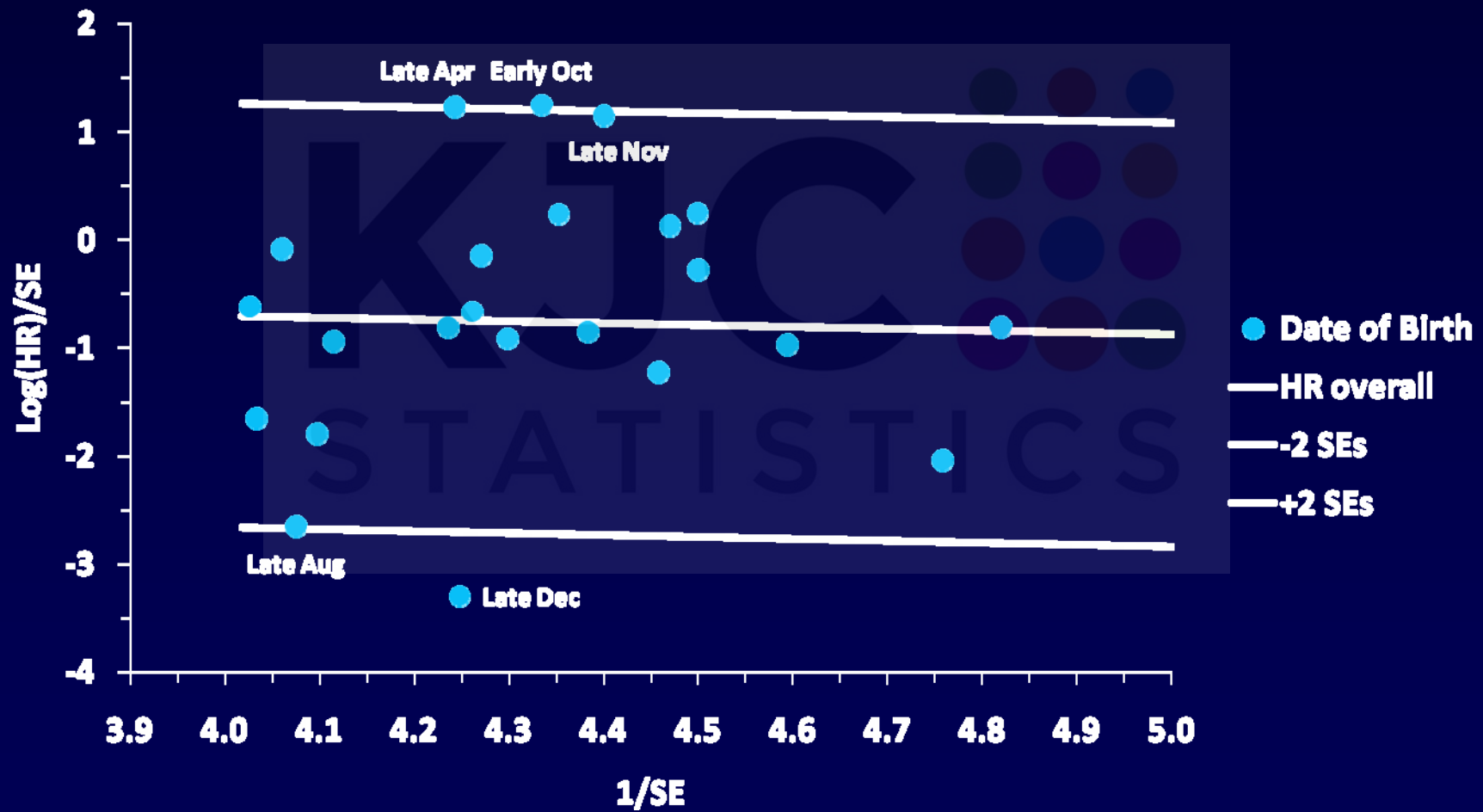
## Probability of Observing At Least 1 Statistically Significant Treatment Interaction By Chance Alone is High

<b>Correlation between tests</b>	<b>Fraction of simulations with at least 1 significant result in 31 tests</b>	<b>Equivalent adjusted p-value to retain overall false-positive error rate at 5%</b>
<b>0</b>	<b>79.1%</b>	<b>0.002</b>
<b>0.5</b>	<b>51.3%</b>	<b>0.004</b>
<b>0.9</b>	<b>17.3%</b>	<b>0.014</b>
<b>0.99</b>	<b>7.9%</b>	<b>0.031</b>
<b>1.0</b>	<b>5.0%</b>	<b>0.050</b>

# PLATO Pattern of effect reversals consistent with what would be expected in a large MRCT

Expected no. countries with HR >1	Actual no. countries with HR >1	Expected no. countries with HR >1.25	Actual no. countries with HR >1.25
12.9	12	6.2	3

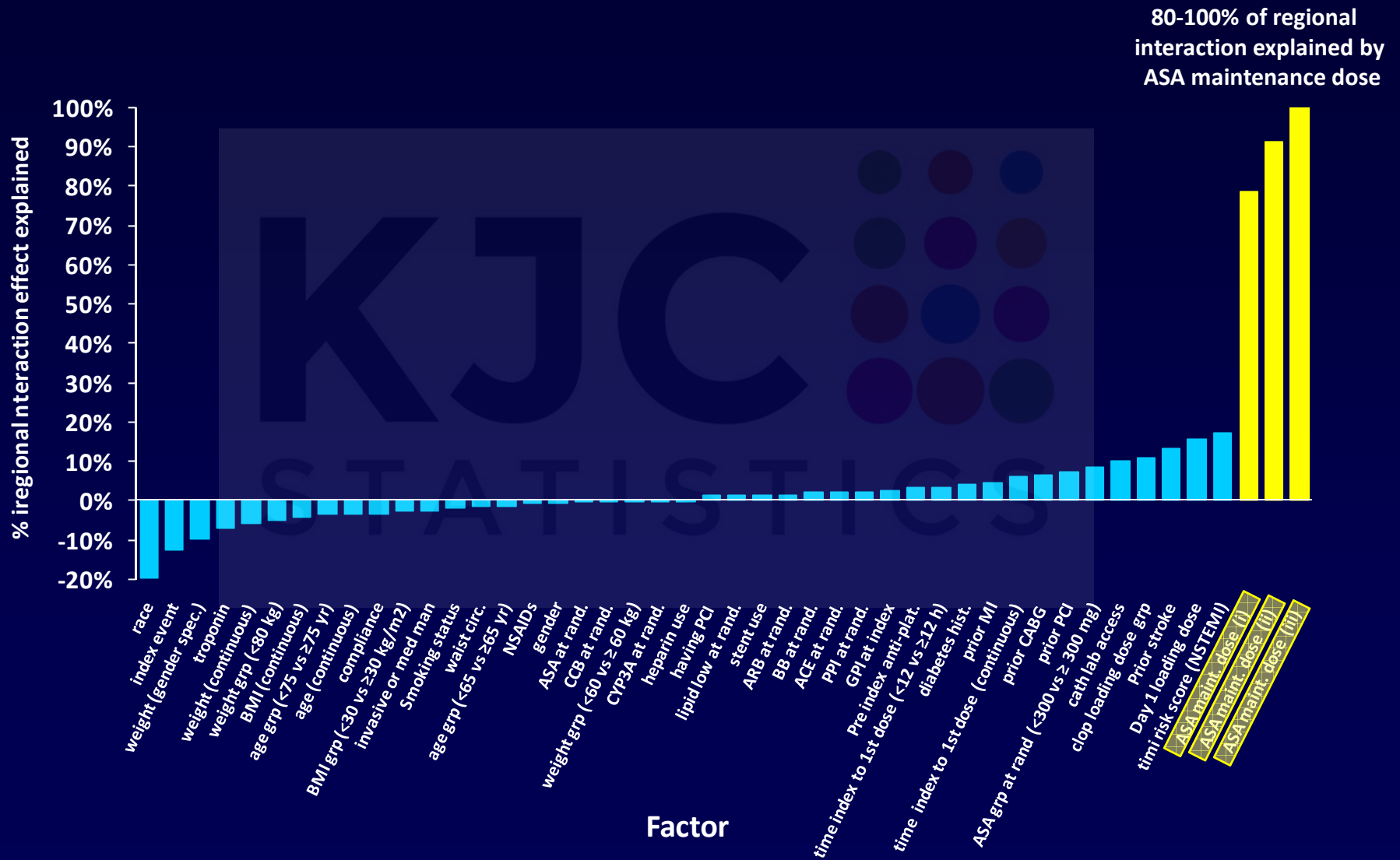
# Analysis by date of birth



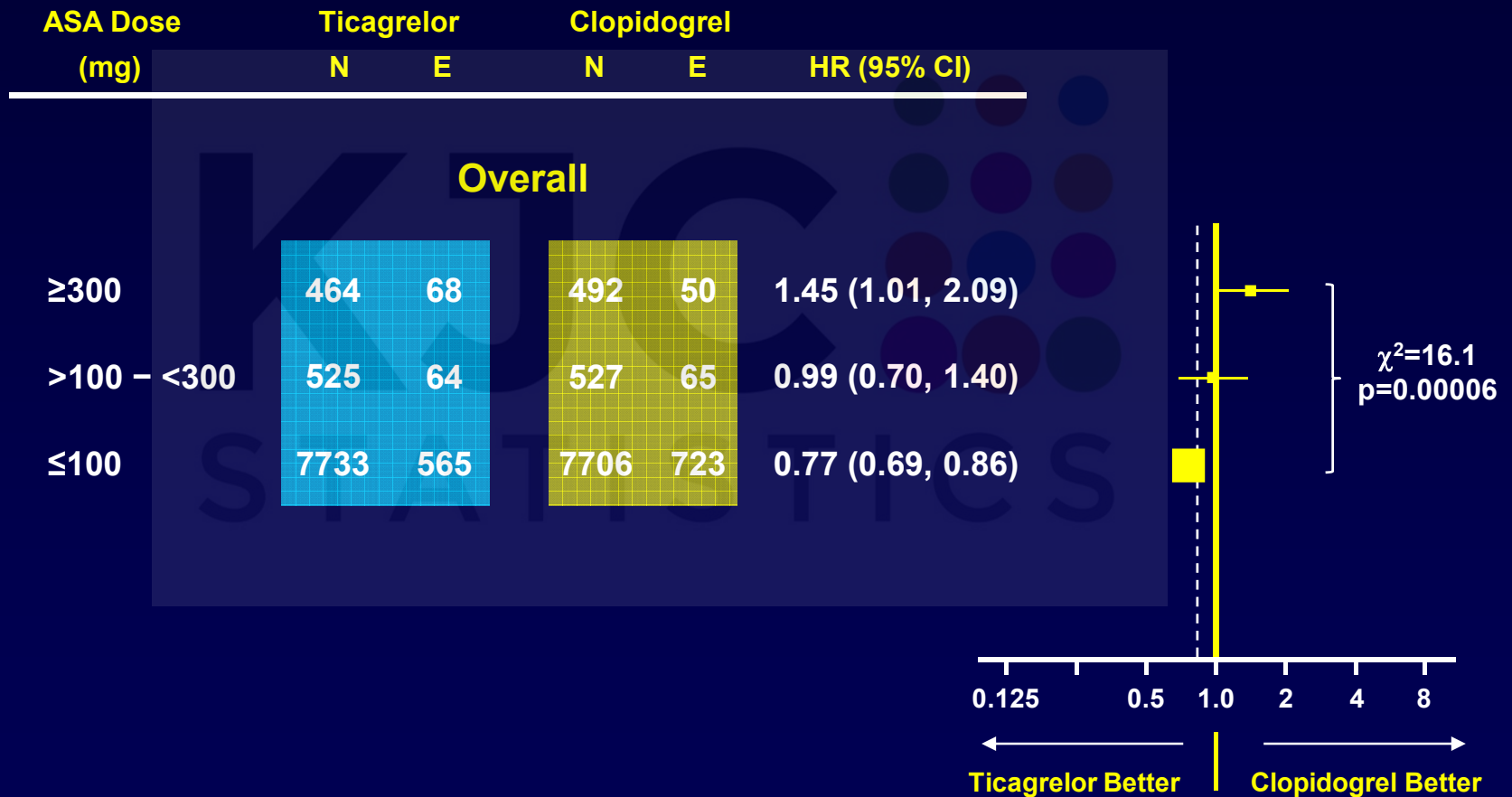
## PLATO: What Kind of Factors or Patient Characteristics Might Explain the US vs Non-US Result?

- To explain a meaningful fraction of the US/non-US interaction, a factor is needed that simultaneously:
  - (i) has a strong qualitative interaction with randomized treatment for the primary endpoint and
  - (ii) is strongly imbalanced between US and non US settings
- Weakly imbalanced prognostic factors will likely not be sufficient to explain the US result
- Visual inspection for imbalances of clinical concern needs to be supported by an objective and statistically rigorous analysis of the data

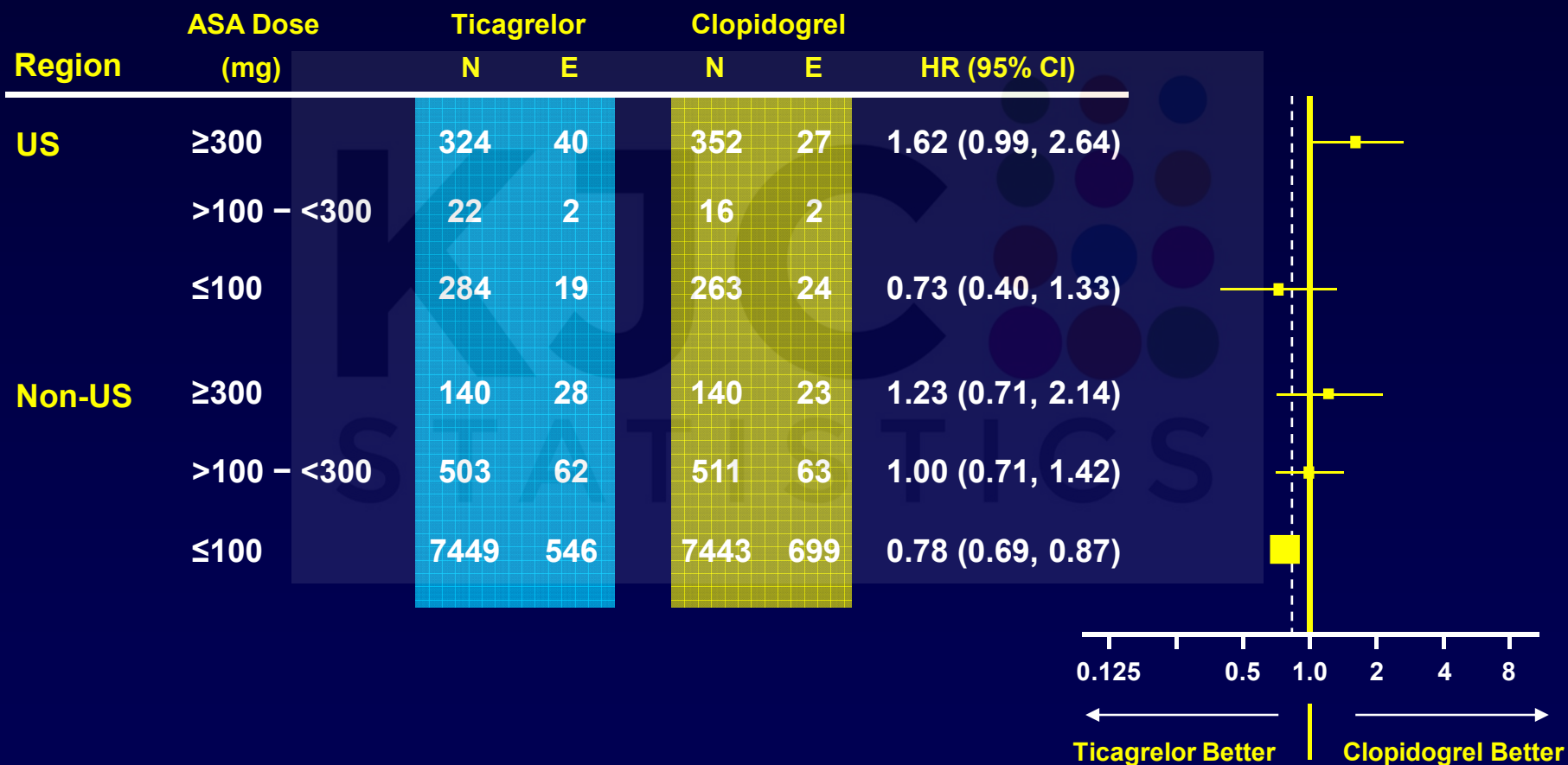
# PLATO: No Factor Potentially Accounts for the Regional Interaction with the Exception of ASA Maintenance Dose During Therapy



# PLATO: The regional interaction is explained by an interaction with ASA maintenance dose



# PLATO: Similar Pattern of Treatment Effects in Relation to ASA Maintenance Dose in US and Non-US



STRICTLY CONFIDENTIAL BRILINTA Regulatory Strategy and Response Document. The preparation for regulatory response requires the discussion of many issues and the preparation for multiple scenarios. This draft document reflects the current thoughts and ideas of multiple individuals involved in the creation or revision of this document as it relates to the subject matter contained herein. It does not necessarily reflect the final view or position of AstraZeneca as it relates to such subject matter and individual views espoused herein may not be based on the full and complete information necessary to make a final determination.

# Summary (1)

- MRCTs are essential in modern day drug development
- Some degree of observed variability is inevitable between regions – in a trial with 80% power and 4 equally sized regions, the chance of at least one reversal is ~30% and with 6 regions the chance is >50%. The corresponding figures are ~20% and 7 regions for a trial with 90% power.
- By definition, MRCTs are not designed, powered or intended to look statistically for true differences between regional effects – all the caveats, biases and pitfalls of subgroup analyses apply.
- Assuming at the planning stage that there is true heterogeneity in the treatment effect by region so ( $\tau^2 > 0$ ) is problematic
  - What value to choose for  $\tau^2$  and on what basis? At what point does the MRCT become meaningless and impossible to interpret clinically due to excessive variability?
  - Sample size inflation will quickly render the trial infeasible.



## Summary (2)

- A random effects analysis is not the answer especially if applying a Follmann adjustment
- More thought needs to be given to the consistency in defining 'region' across trials – “Brazil, Chile, Argentina and New Zealand”
- The optimum allocation of N across regions / countries and the criteria for 'consistency' are two sides of the same coin.
- For multiple regions, simplest, but least realistic, criteria is to require all regional effects  $> 0$ .
- For reference vs non-reference regions, most pragmatic criteria is to require 3 regions, is to require the point estimate for the reference region to be no less than the lower  $100(1-2\alpha)\%$  CI estimate for the treatment effect in the total trial population.
  - provides some reassurance the treatment effect in the reference region is not worryingly less than the treatment effect seen in the total trial population.

## Summary (3)

- Graphical methods are very informative when examining regional effects – Galbraith, Q-Q and forest plots are essential.
- Decomposition of the overall  $\chi^2$  into a weighted sum of  $t^2$  statistics for  $i^{\text{th}}$  region vs all others is very informative and should be done routinely.
- If there is truly believed to be a region with a different result, extensive evaluation is required to assess play of chance, and to look for possible confounders in medical practice, medical care, quality of trial conduct, cultural, biologic and genetic factors that may be explanatory.
- In the end “Thus the best estimate of the treatment effect ... for any subgroup is the estimate of the ... [treatment effect] for the overall trial” Wedel and DeMets