

1 Sample size calculations for randomised controlled trials and for prediction models

2

3 Francesca Fiorentino

4 Department of Surgery and Cancer, Imperial College London, London UK

5

6 Tom Treasure\*

7 Clinical Operational Research Unit, University College London

8

9

10 Corresponding author

11 Tom Treasure: [tom.treasure@gmail.com](mailto:tom.treasure@gmail.com)

12

13

14

15

16

17

18 There was no funding for this work.

19 The authors have no conflicts of interest to declare.

20

21 Abstract

22

23 The two study protocols are published in this issue Colorectal Disease: FALCON, a  
24 multicentre randomised controlled trial of strategies to reduce surgical site infection, and  
25 AFAR, a predictive model of atrial fibrillation after colonic resection. Both are exemplars of  
26 excellent research design that surgeon researchers should seek to emulate. Trial statisticians  
27 were involved at an early stage and the protocols have been through several rounds of peer  
28 review by trial methodologists, prior to being funded by the National Institute for Health  
29 Research (NIHR). In this article we address the important question of sample size  
30 calculations and how they should be approached for these very different forms of study.

31

32

33

34 Main Text

35

36 Most surgical procedures came into practice without randomised trials because, against a  
37 well-known experience of clinical outcomes over many years, an appropriate and well  
38 conducted operation was seen to make a dramatic and lasting difference. For example,  
39 Thomas's splint only had to be seen in use for injured farmers in north Wales, and then  
40 soldiers in the 1914-18 war, to become universally adopted. The relief of pain in the hours  
41 and days after injury was evident, followed by recovery to walk on legs of matching length,  
42 with both feet pointing forward. To generalise that process of deduction, the features that  
43 indicate that an RCT is not needed are a close temporal and mechanistic relationship between  
44 the intervention and the effect, resulting in a large and sustained benefit.(1) The Thomas's  
45 splint became the standard initial treatment, applicable to the large majority of patients with  
46 femoral fracture.

47

48 In contrast, lung metastasectomy is carried out in fewer than one in thirty of the patients who  
49 have lung metastases.(2) The outcome of importance is survival. For lung metastasectomy,  
50 results are usually given as survival rate, usually at an interval of five-years, but there are too  
51 many factors and uncertainties to conclude that metastasectomy has a survival benefit by  
52 observation alone.(3)

53

54 **Calculating the sample size for a randomised controlled trial**

55 It is wasteful of time and effort to embark on a study that is not large enough to provide a  
56 conclusive answer, or so large as to be wasteful of effort and resources.(4) To calculate a  
57 sample size, the statistician needs to know what is (a) the outcome of importance, (b) the  
58 outcome measure and (c) the clinically meaningful effect size.

59

60 For lung metastasectomy, survival beyond five years was the only outcome reported in the 51  
61 follow-up studies found in a systematic review (5) so for our first illustration (a) survival is  
62 the outcome of importance. Survival of ~40% at five years has been consistently reported and  
63 was confirmed in a meta-analysis including 2925 patients. (6) For the illustration we will  
64 identify the survival rate at 5-years to be the outcome measure (b). The effect size depends on  
65 what would be the survival without metastasectomy. The US Society of Thoracic Surgeons  
66 based its recommendations on a consensus assumption of zero survival, but for this  
67 illustration will use the more cautious "worse than 5%" suggested by the authors of the meta-  
68 analysis. Then (c) is the absolute difference between 40% and 5%, the effect size of 35%.

69

70 The surgeons need to agree with statistician the value of *alpha*—the probability of a false  
71 positive—usually set at 5% and hence the familiar  $P < 0.05$ . The value of *beta*—the probability  
72 of a false negative—is usually set at 20% or more cautiously 10%. Power is  $1 - \beta$  so in  
73 percentage terms these are expressed as 80% or 90%, that is the power to avoid a false  
74 negative. Given these estimates a statistician can generate Table 1. This is for 1:1  
75 randomisation and shows that 44 patients (22 in each arm) would provide 80% power for a  
76 two-sample proportion test. There are likely to be patients lost to follow up, so the target  
77 recruitment might be set at 50.

78  
79 In cancer trials it is usual to use time to death (overall survival) or cancer progression  
80 (progression free survival) for (b) the outcome measure. The statistical test used for the  
81 sample size calculation is the two-sample comparison of survivor functions (log-rank test).  
82 The same assumptions can be used to do the calculation, but the statistical method takes into  
83 account the time of the event, death. It captures more information than a simple count of 5-  
84 year survivors, so it requires commensurately fewer patients. Using the log rank test, the  
85 statistician can produce Table 2. Randomisation is still 1:1 and shows that 36 patients (18 in  
86 each arm) would provide 80% power with a two-sample survivor function test. A total of 42  
87 patients would allow for loss to follow up.

88  
89 In the discussion between the investigators and the statistician, all should be alert to the  
90 possibility of “back calculation”. The surgeons know the number of patients available for  
91 recruitment and can tweak the effect size to give an achievable number of randomised  
92 patients. In the case of lung metastasectomy the consensus assumption of zero survival(7) had  
93 for years ruled out the possibility of randomisation at all; there was no prospect of equipoise.  
94 Also, it conveniently attributed all the credit for survival to the effect of the operation and  
95 trumps any likely effect from chemotherapy.

96  
97 Tables 1 and 2 illustrate the principle, but it is not how the conversation with the statistician  
98 went in the case of the PulMiCC trial. The investigators had reason to believe that patients  
99 eligible for metastasectomy had better survival than was widely assumed. This came from a  
100 comparative study in 1980(8) and a modelling study on cancer registry data in 2006.(9) Both  
101 suggested the possibility that metastasectomy makes a much smaller difference to survival  
102 than assumed. Knowing that, the statistician asked what was the smallest clinically  
103 meaningful difference in the five-year survival that might justify lung metastasectomy. A  
104 difference from 40% survival in the treated down to 30% survival in the control was  
105 suggested (10% difference). Table 3 shows the calculation using a two-sample survivor  
106 function test.

107  
108 As we said, the actual sample size calculation may be much more complicated. In fact, the  
109 PulMiCC trial was powered for *non-inferiority* of leaving the metastases unresected using  
110 time to event analysis.(10) With this smaller difference (40% and 30%) the numbers needed  
111 to power the study were commensurately higher, and in the event not achievable due to the  
112 tenacity with which cancer teams held on to the near zero assumption and its implications.  
113 (11, 12) It is also important to remember that for the sample size calculation it was important  
114 to be realistic at the planning stage, The assumptions are replaced by findings once the data  
115 are in, and the prior power calculation plays no part in the analysis or interpretation of the  
116 results.(13)

117  
118 It may be important to not rely on randomisation, but to ensure that there is a balance in  
119 prognostic factors between the randomized groups, particularly if these factors might create

120 differences of a magnitude that compete with the treatment effect (confounding factors). For  
121 example, obesity in studies of surgical site infection(14) which might be relevant in  
122 FALCON.(15) In the case of the PulMiCC trial the unfavourable features were more than one  
123 metastasis, liver involvement, carcinoembryonic antigen elevation and shorter interval since  
124 the primary resection. In large drug trials this process is done by stratification but in trials of  
125 limited size an alternative is *minimisation* which adjusts the probability of a patient being  
126 assigned to one or other arm in order to achieve balance between the groups in the known  
127 factors, relying on randomisation to balance the unknown confounders.(16, 17) It is essential  
128 that this is done by a strict algorithm out of sight of anyone involved in the trial.

## 131 Prediction models

132  
133 Prediction models are used for investigating patient outcomes in relation to patient and  
134 disease characteristics. They may be of use in surgical practice and we give three examples.

- 135  
136 1. In the AFAR study(18) the adverse outcome to be “predicted” is the onset of new  
137 atrial fibrillation during the recovery period. Patients in the stratum more likely to  
138 have this problem can then have further planned screening or prophylactic  
139 approaches. The model is intended to target more costly and labour intensive methods  
140 to where they will achieve the greatest benefit for patients.
- 141 2. A predictive model has been developed to risk adjust postoperative mortality among  
142 patients having of colorectal cancer.(19) It allows fair comparisons to be made  
143 between hospitals, clinical teams and individual surgeons. Implementation of public  
144 reporting in 2013 was followed by a fall in the observed surgical mortality. The model  
145 allowed this to be interpreted as a real reduction in mortality without risk avoidance.  
146 (20)
- 147 3. A third example is to select patients for surgery by gaining insights into their  
148 likelihood of death or survival after surgery. We will return to an unsatisfactory  
149 example in the development of a model with this purpose as a cautionary tale.(21)

150  
151 A standard approach is to use a “training” dataset and the model is then tested with a separate  
152 “validation” dataset which has been held back for the purpose. Following the same principle  
153 as the sample size calculation the statistician must be provide with the best available data,  
154 informed estimates of as yet unquantified factors and what outcome would be useful. The  
155 outcome can be a continuous scale, categorical or estimated survival (time to event).

156  
157 The model developed by Walker, Finan and van der Meulen(22) used internal validation and  
158 is more sophisticated than can be described here but it illustrates the power of a collaborative  
159 effort with data available on 62,314 patients in the National Bowel Cancer Audit and  
160 collaboration with very highly skilled data analysts. Eight risk factors were included and  
161 mortality was counted up to 90 days. This captures 50% more deaths, virtually all having a  
162 relationship to treatment. The methods of “imputation” for missing data (and missingness is  
163 inevitable) and validation were at a high level of expertise.

164  
165 The tables are set up to illustrate the fewest of counts that might allow for a valid model. In  
166 the case of survival (time to event outcome), the overall event rate and the mean follow-up  
167 time need to be known. In the case of binary outcome, the outcome proportion expected  
168 within the model development dataset, based on previous evidence.

169 Tables 4 and 5 give examples of sample size for prediction models of binary and survival  
170 outcomes prediction models.

171

172 For less common disease or particular circumstances there may still be a desire to create  
173 models to inform practice. A recent published example is of a scoring system to select  
174 patients more likely to “benefit” from lung metastasectomy for sarcoma included 135  
175 patients.(21) The scoring system has three parameters giving scores of 0-3. What can be  
176 lauded in the report is that the authors provide the data. The figure is taken from their paper.  
177 The well-used caution “correlation does not mean causation” can be applied. The more  
178 important value of  $r^2$  is 0.144 indicating that the scores contribute very little, <15%, to the  
179 prognosis. It is clear from the graphical display that the scores really do not discriminate  
180 usefully between lengths of survival.

181

182 Sarcoma has a predilection for metastasising to the lungs and these patients are often young  
183 so there is pressure to do something, anything, to help. The two longest survivors at over 12  
184 years who scored 2/3 on the scale will be pointed out repeatedly on clinic visits, generating  
185 confirmation bias. What will not be recalled is the harm done by operations of unproven  
186 effectiveness on patients the larger number of patients not coming back to clinic.

187

188

189 Table 1

Power	Alpha	Effect size	N
80%	0.05	0.35	44
90%	0.05	0.35	56

190

191 In the case of an effect size of 35% (0.35) and using a two-sample proportions test (Pearson's chi-squared test)  
192 the variation in power and sample size can be seen, for the first scenario with assumed 35% survival gain from  
193 metastasectomy over control.

194

195 Table2

alpha	power	N	Expected events	Hazard Ratio	Survival Metastasectomy	Survival no operation
80%	0.5	36	28	3.269	40%	5%
90%	0.5	48	38	3.269	40%	5%

196

197 A sample size calculation for the two-sample comparison of survivor functions (log-rank test) using the same  
198 assumptions.

199

200 Table 3

Power	Alpha	N	Expected events	Hazard Ratio	Survival Metastasectomy	Survival no operation
80%	0.5	656	427	1.314	0.4	0.3
90%	0.5	880	571	1.314	0.4	0.3

201

202 A sample size calculation for the two-sample comparison of survivor functions (log-rank test) raising the control  
203 estimate from 5% to 30%.

204

205

206 Table 4

Predictor parameters	Outcome prevalence	Minimum Sample Size	Number of events
10	10%	348	35
20	10%	695	70
30	10%	1042	105
10	40%	369	148
20	40%	519	208
30	40%	778	312

207

208

209

210

211 Example for a binary outcome where the expected outcome proportion is 10% or 40% with model parameters  
212 10, 20 and 30 in Table 4.

213

214

215

216 Table 5

Predictor parameters	Overall event rate	Minimum Sample Size	Number of outcome events
10	6.5%	1715	231
20	6.5%	3429	462
30	6.5%	5143	692

217

218

219

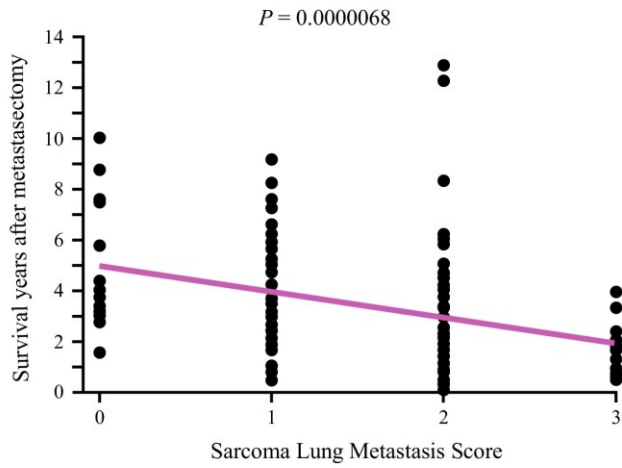
220

221

222 Example for a survival outcome where the mean follow-up is 2 years, the overall event rate is 0.065 and the  
223 time for model prediction is 2 years with model parameters 10, 20 and 30 in Table 5.

224

225



**FIG. 4** Relationship between the survival period after lung metastasectomy and the Sarcoma Lung Metastasis Score. Pearson's product-moment correlation coefficient was calculated between the Sarcoma Lung Metastasis Score and survival time after lung metastasectomy, and a significant correlation was observed between these two factors ( $r = -0.38$ ,  $p = 0.000068$ )

226  
 227  
 228  
 229



230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279

1. Glasziou P, Chalmers I, Rawlins M, McCulloch P. When are randomised trials unnecessary? Picking signal from noise. *BMJ*. 2007;334(7589):349-51.
2. Fenton H, Finan PJ, Milton R, Shackcloth M, Taylor J, Treasure T, et al. National variation in pulmonary metastasectomy for colorectal cancer. *Colorectal Disease*. 2020;Accepted for Publication.
3. Milosevic M, Edwards J, Tsang D, Dunning J, Shackcloth M, Batchelor T, et al. Pulmonary Metastasectomy in Colorectal Cancer: updated analysis of 93 randomized patients - control survival is much better than previously assumed. *Colorectal Dis*. 2020;22(10):1314-24.
4. Glasziou P, Chalmers I. Research waste is still a scandal-an essay by Paul Glasziou and Iain Chalmers. *BMJ*. 2018;363:k4645.
5. Fiorentino F, Hunt I, Teoh K, Treasure T, Utley M. Pulmonary metastasectomy in colorectal cancer: a systematic review and quantitative synthesis. *J R Soc Med*. 2010;103(2):60-6.
6. Gonzalez M, Poncet A, Combescure C, Robert J, Ris HB, Gervaz P. Risk factors for survival after lung metastasectomy in colorectal cancer patients: a systematic review and meta-analysis. *Ann Surg Oncol*. 2013;20(2):572-9.
7. Handy JR, Bremner RM, Crocenzi TS, Detterbeck FC, Fernando HC, Fidas PM, et al. Expert Consensus Document on Pulmonary Metastasectomy. *Ann Thorac Surg*. 2019;107(2):631-49.
8. Aberg T, Malmberg KA, Nilsson B, Nou E. The effect of metastasectomy: fact or fiction? *Ann Thorac Surg*. 1980;30(4):378-84.
9. Utley M, Treasure T, Linklater K, Moller H. Better out than in? The resection of pulmonary metastases from colorectal tumours. In: Xie X, Lorca F, Marcon E, editors. *Operations Research for Health Care Engineering: Proceedings of the 33<sup>rd</sup> International Conference on Operational Research Applied to Health Services*. Saint-Etienne: Publications de l'Universitaire de Saint-Etienne; 2008. p. 493-500.
10. Milosevic M, Edwards J, Tsang D, Dunning J, Shackcloth M, Batchelor T, et al. Pulmonary Metastasectomy in Colorectal Cancer (PulMiCC): Updated analysis of 93 randomised patients - control survival is much better than previously assumed. *Colorectal Dis*. 2020;<https://onlinelibrary.wiley.com/doi/full/10.1111/codi.15113>.
11. Treasure T, Farewell V, Macbeth F, Monson K, Williams NR, Brew-Graves C, et al. Pulmonary Metastasectomy versus Continued Active Monitoring in Colorectal Cancer (PulMiCC): a multicentre randomised clinical trial. *Trials*. 2019;20(1):718.
12. Macbeth F, Fallowfield L. The myth of pulmonary metastasectomy. *Br J Cancer*. 2020;123(4):499-500.
13. Cox DR. *Planning of Experiments*: Wiley; 1958.
14. Wilson AP, Livesey SA, Treasure T, Gruneberg RN, Sturridge MF. Factors predisposing to wound infection in cardiac surgery. A prospective study of 517 patients. *Eur J Cardiothorac Surg*. 1987;1(3):158-64.
15. Nepogodiev D, Bahngu A. Pragmatic multicentre factorial randomised controlled trial testing measures to reduce surgical site infection in low- and middle-income countries: study protocol of the FALCON trial. *Colorectal Disease*. 2021;IN PRESS.
16. Treasure T, MacRae KD. Minimisation: the platinum standard for trials? *BMJ*. 1998;317(7155):362-3.
17. Altman DG, Bland JM. Treatment allocation by minimisation. *BMJ*. 2005;330(7495):843.
18. Lee MJ, Hawkins DJ, Bradburn MJ, Lee J, Brown SR, Wilson MJ. Atrial Fibrillation After Resection (AFAR) A progress III Study. *Colorectal Disease*. 2021;IN PRESS.
19. Walker K, Finan PJ, van der Meulen JH. Model for risk adjustment of postoperative mortality in patients with colorectal cancer. *Br J Surg*. 2015;102(3):269-80.

- 280 20. Vallance AE, Fearnhead NS, Kuryba A, Hill J, Maxwell-Armstrong C, Braun M, et al. Effect of  
281 public reporting of surgeons' outcomes on patient selection, "gaming," and mortality in colorectal  
282 cancer surgery in England: population based cohort study. *BMJ*. 2018;361:k1581.
- 283 21. Yamamoto H, Yamamoto H, Soh J. A Simple Prognostic Benefit Scoring System for Sarcoma  
284 Patients with Pulmonary Metastases: Sarcoma Lung Metastasis Score. *Annals of Surgical Oncology*.  
285 2020.
- 286 22. Vallance AE, Van Der Meulen J, Kuryba A, Braun M, Jayne DG, Hill J, et al. Socioeconomic  
287 differences in selection for liver resection in metastatic colorectal cancer and the impact on survival.  
288 *Eur J Surg Oncol*. 2018;44(10):1588-94.

289