

1 **External model validation of binary clinical risk prediction models in**
2 **cardiovascular and thoracic surgery**

3 Graeme L. Hickey¹, Eugene H. Blackstone²

4 ¹ University of Liverpool, Department of Biostatistics, Waterhouse Building (Block F),
5 1-5 Brownlow Street, Liverpool, L69 3GL, United Kingdom

6 ² Cleveland Clinic, Heart and Vascular Institute Clinical Investigations, 9500 Euclid
7 Avenue, JJ4, Cleveland, OH 44195, United States

8

9 The authors have no conflict of interest with material in this submission.

10

11 **Words: ~ 2000**

12

13 **Corresponding Author:**

14 Dr Graeme Hickey

15 University of Liverpool, Department of Biostatistics

16 Waterhouse Building (Block F), 1-5 Brownlow Street

17 Liverpool, L69 3GL, UK

18 **Email:** graeme.hickey@liverpool.ac.uk

19 **Tel:** +44 (0)151 794 9737

20

21

CENTRAL PICTURE

22

23

[Use Figure (central).pdf]

24

25
26
27
28
29

CENTRAL MESSAGE

[Characters + spaces = 117; Limit = 200]

30 External validation of binary clinical risk-prediction models is vital. We provide
31 strategies for accomplishing this.

32
33

34 **INTRODUCTION**

35 Clinical risk-prediction models (CRPMs, also known as prognostic models or
36 risk score models) serve an important role in healthcare,¹ particularly for binary
37 adverse events (in-hospital, 30-day, or operative mortality) after cardiac, thoracic,
38 and vascular surgery. These models may be applied to 3 different objectives: 1) to
39 assess patient risk, which surgeons and patients can then factor in to healthcare
40 decisions; 2) to stratify risk, both for clinical decision-making and inclusion criteria in
41 a controlled randomized trial,² and 3) to assess and compare healthcare outcomes
42 among providers (benchmarking). The comparison of observed with expected
43 outcomes, accounting for statistical uncertainty, can identify underperforming
44 healthcare providers for quality improvement interventions.³

45 The wide-ranging importance of CRPMs in the cardiovascular specialty
46 means that stakeholders must have confidence in them. A poorly performing model
47 can lead to suboptimal decision-making, misinformed patients, false reassurance of
48 a healthcare provider's performance, or false stigmatization of the provider.
49 Confidence is established by validating the model.⁴

50 Model validation can be internal, temporal, or external. Internal model
51 validation is one element of CRPM development, usually published alongside the
52 model to confirm the model performs well for the training data. External validation,
53 which evaluates the generalizability (or transportability) of the model to other groups
54 of patients, is fundamental to demonstrating a model is appropriate for adoption in
55 clinical practice.⁴ In cardiovascular and thoracic surgery, the majority of CRPMs
56 encountered will predict binary outcomes, which were created using multivariable
57 regression techniques, in particular logistic regression. Therefore, we focus our
58 discussion to this area. However, the general principles and need for external

59 validation apply to other outcome types and models, e.g. time-to-event data,^{5,6} as
60 well as to non-regression techniques, e.g. machine learning approaches.⁷

61

62 **MODEL PERFORMANCE CONCEPTS**

63 Performance of CRPMs is typically based on assessing two important
64 features: calibration and discrimination.⁶

65 **Calibration** is the accuracy of the model for predicting events relative to
66 observed events in groups of patients. For example, if the mean predicted event
67 occurrence is 5% in a patient group, but the observed event occurrence is 10%, then
68 we conclude the model is not well calibrated because it underpredicts.

69 **Discrimination** is the ability of a model to distinguish between patients who
70 experienced the event and those who did not. Discrimination is measured using the
71 area under the receiver-operating-characteristic curve (AUROC), also referred to as
72 the concordance (c-)statistic or c-index.⁵ This value has a meaningful interpretation.
73 If we randomly select 2 patients, 1 who experienced the event and 1 who did not,
74 then the AUROC is equivalent to the probability that the risk score attributed to the
75 former is greater than that attributed to the latter. An AUROC of 1 indicates perfect
76 classification; a value of 0.5 is equivalent to tossing a fair coin.

77 Other aspects of performance assessment include clinical usefulness,
78 impact,⁸ and overall performance measures such as the Brier score⁹ and
79 concordance index, particularly for time-related events.

80

81 **DESIGNING AND REPORTING AN EXTERNAL VALIDATION**

82 When designing a validation study, thought must be given to several key
83 elements.

84 **Selection of patients.** The selection of patients used to externally validate a
85 CRPM might differ from those used to develop the model. These differences might
86 be temporal or geographical, or related to clinical setting, inclusion or exclusion
87 criteria, definitions, diagnostic techniques, or inherent baseline case-mix differences
88 between the two populations. It is important to highlight any differences that might
89 affect model transportability between the validation and original study sample,
90 particularly with validation of general all-surgery models (e.g. the EuroSCORE)
91 within procedural¹⁰ or operative subgroups.¹¹

92 **Risk factor data.** It goes without saying that calculating a risk score requires
93 access to all variables that comprise the risk score. One potential issue is conflict in
94 variable definitions. For example, a registry that only collects binary data on whether
95 pulmonary artery (PA) systolic pressure is >60 mmHg (a risk factor in the logistic
96 EuroSCORE model) would not be able to compute the EuroSCORE II risk score,
97 which includes model coefficients for PA systolic pressures of 31 – 55 mmHg and
98 >55 mmHg. This is primarily an issue for retrospective validation studies, as clinical
99 registries can be updated to capture contemporary risk-score data.

100 **Missing data.** One cannot calculate a risk score without access to data for
101 variables that comprise the CRPM. If a model contains a risk factor such as
102 preoperative serum creatinine, but these data are sparsely available in the dataset,
103 then in many cases the risk score cannot be calculated. Case-complete analyses—
104 those that delete subjects with missing data for required variables—might lead to
105 bias if those subjects are not representative of the whole population.¹² In certain
106 cases, reasonable estimates and assumptions can be made based on clinical
107 expertise or additional information in the dataset. For example, a number of variables
108 in Society of Thoracic Surgeons (STS) risk models have coefficients set to 0 for

109 some variables in some models; if one is validating such a model, missing data for
110 such a variable is of no consequence. Alternatively, statistical imputation or subset
111 analysis techniques might be applied to compensate.^{13,14} If a validation study
112 specifically excludes certain groups of patients (for example, emergency surgery,
113 reoperations, or endocarditis), imputation of 0 is an accurate and appropriate
114 substitution, but the validation is only partial. In any case, it is always necessary to
115 summarize the frequency of missing data and present methods for managing it and
116 its assumptions.

117 **Sample size.** Considerations regarding sample size should not be limited to
118 randomized control trials. Single-center validation studies will often have a limited
119 pool of subjects, especially for subgroup analyses, and increasing the sample size
120 will require widening the study period, which could come at a price (see comment on
121 calibration drift below). When designing a study, sample size (number of subjects)
122 alone is not enough; one must also consider effective sample size (number of
123 events). Relatively little attention has been given to this matter, but some studies
124 have recommended a minimum of 100 events and 100 non-events for validation
125 studies, and in certain applications, larger effective sizes will be required to obtain
126 adequate power.^{15,16}

127 **Outcome definitions.** Many well-known CRPMs in cardiac surgery predict
128 early or operative mortality, including the logistic EuroSCORE¹⁷ and STS Cardiac
129 Surgery Risk Models.^{18–20} Operative mortality is generally accepted to mean death
130 within 30 days (or later if the patient has not been discharged within 30 days).²¹
131 However, other definitions of mortality exist, such as in-hospital mortality.²² Two
132 large databases reported operative mortality to be 4.63% and 3.57%, compared with
133 in-hospital mortality of 4.02% and 2.94%, respectively.^{23,24} In both cases, in-hospital

134 mortality was approximately 0.6% lower. In-hospital mortality is generally easier to
135 robustly measure, whereas 30-day mortality requires post-discharge follow-up for
136 most patients.²⁵ Therefore, it is common to see models validated against in-hospital
137 mortality. In this example, we would expect the model to over-predict mortality
138 relative to the observed data. It is reasonable to assess the model performance for
139 this similar endpoint; however, this subtlety should be borne in mind when designing
140 a study, particularly if the objective of the study is to compare models that have
141 different outcome definitions. Similar considerations apply to cases where the
142 definition of a major postoperative complication used for model development differs
143 from that in the validation dataset.

144 **Large study windows.** One simple way to increase sample size in a
145 validation study is to widen the study window. However, validation of a CRPM over a
146 substantially wide period can introduce a number of complexities. One potential
147 issue is calibration drift.^{26,27} Multiple studies demonstrated that the ratio of observed
148 mortality to mean logistic EuroSCORE was decreasing with time. Changing risk
149 profiles, other variables influencing mortality, and changes in the association of risk
150 factors with outcome can all contribute to this phenomenon. This prompted the
151 introduction of the EuroSCORE II model²³ and the series of contemporary STS
152 models.^{18–20} Researchers should be aware of this, particularly when validating
153 cardiac surgery CRPMs.

154 **TRIPOD statement.** In recent years, reporting of biomedical research has
155 been improved with guidelines such as the CONSORT statement²⁸ for randomized
156 trials and the PRISMA statement²⁹ for systematic reviews and meta-analyses.
157 Prompted by evidence of poor quality reporting in the CRPM literature, the recent
158 TRIPOD statement describes reporting guidelines for studies developing, validating,

159 or updating a prediction model.³⁰ We strongly encourage researchers to follow these
160 guidelines and make use of the checklist for validating models. Examples of good
161 practice and additional details have been previously published.³¹

162

163 **METHODS FOR ASSESSING CALIBRATION**

164 ***Hosmer-Lemeshow test.*** The Hosmer-Lemeshow test is a frequently
165 reported statistical test for assessing calibration in CRPMs. However, it has a
166 number of drawbacks.^{31–35} First, it is not easily interpreted; that is, it does not provide
167 a measure of the magnitude of any miscalibration. Second, for slight deviations in
168 calibration, the test is sensitive to sample size. Third, the classical version of the test
169 is dependent on arbitrary groupings of patients. In some cases, the Hosmer-
170 Lemeshow test remains a useful adjunct statistic, but should only be included as part
171 of a more comprehensive assessment. Typically, the Hosmer-Lemeshow test refers
172 to a test based on 10 groups composed by deciles of risk. However, authors should
173 be aware that there are variations on the test with regard to groupings (quantiles vs.
174 fixed cut-points), number of groups (g), degrees of freedom of the chi-squared
175 statistic ($g-2$ for internal vs. g for external validation), and software
176 implementations.^{35,36} While g is typically selected to be 10, one must ensure the cell
177 counts are sufficient to justify the distributional approximation. Including a table of
178 observed and expected events by binning group provides a useful summary, and
179 allows for inspection of each term for fit, as recommended by Hosmer and
180 Lemeshow (p. 188).³⁶

181 ***Calibration plot.*** If a standard Hosmer-Lemeshow test is performed, then a
182 simple graph—the calibration plot—is a straightforward next step (**Figure**).⁴ Within
183 each of the g groups, observed events are plotted against expected events. If the

184 model is well calibrated, then these points should be close to the 45° line. The
185 calibration plot can be augmented by overlaying a non-parametric smoothing curve
186 (e.g. loess) through the observed and predicted data³⁷ or a calibration curve.³⁸
187 Contrary to the Hosmer-Lemeshow test and basic calibration plot, these additional
188 fits are not dependent on arbitrary groupings.

189 **Calibration curves.** Cox's calibration regression fits a logistic regression
190 between the observed event and the log-odds transformed predicted values.³⁹ A
191 perfectly calibrated CRPM (deriving from a logistic regression model) yields an
192 intercept = 0 and a slope = 1. These fitted regression models can be superimposed
193 onto a calibration plot, giving an alternative graphical description of the
194 miscalibration. As well as quantifying the degree of miscalibration, one can also
195 simultaneously test whether the estimated parameters reject the null hypothesis of
196 calibration. There are other related null hypotheses that can be tested for assessing
197 calibration also (p. 274).⁶

198 **Other tests.** The Hosmer-Lemeshow is ubiquitous in biomedical CRPM
199 literature. However, researchers can take advantage of a wide variety of statistical
200 tests to assess model validation, such as the aforementioned calibration curve
201 test(s), the Spiegelhalter Z-test,⁴⁰ and methods proposed by Stallard.⁴¹ Most can be
202 calculated using routine software packages.^{6,38} There is no omnibus test of
203 calibration; each approach has different merits and limitations. Therefore, it is
204 important that researchers employ a broad repertoire of methods to address the
205 study questions.

206

207 **MODEL UPDATING**

208 A natural extension to the validation of a CRPM is the concept of updating an
209 existing model. This might involve exploring whether a new biomarker improves a
210 model (e.g. using net reclassification improvement measures⁴²), recalibrating a
211 model,⁴³ and, more recently, assessing whether multiple models can be combined to
212 provide a more accurate prediction (e.g. meta-models and model averaging).⁴⁴ This
213 expanding research area is especially important in an era of personalized
214 medicine.⁴⁵

215

216 **CONCLUSIONS**

217 External validation of CRPMs is necessary to demonstrate their predictive
218 accuracy. Available models have likely been validated internally; however, using
219 them in different settings, locations, and populations can result in relatively poor
220 performance. CRPMs that have been overfitted during development will also often
221 fail to generalise to the external validation sample. Calibration and discrimination
222 must be measured in order to establish validity. There are multiple statistical
223 approaches available to interrogate the calibration, with it being widely accepted that
224 the ubiquitous Hosmer-Lemeshow test has limited utility. Execution of a rigorous
225 CRPM validation study rests in proper study design, application of suitable statistical
226 methods, and transparent reporting.

227

228

REFERENCES

- 229 1. Moons KGM, Royston P, Vergouwe Y, Grobbee DE, Altman DG. Prognosis
230 and prognostic research: what, why, and how? *BMJ*. 2009;338(b375):1317-
231 1320.
- 232 2. Leon MB, Smith CR, Mack MJ, Miller DC, Moses JW, Svensson LG, et al.
233 Transcatheter aortic-valve implantation for aortic stenosis in patients who
234 cannot undergo surgery. *N Engl J Med*. 2010;363(17):1597-1607.
- 235 3. Bridgewater B, Hickey GL, Cooper G, Deanfield J, Roxburgh JC. Publishing
236 cardiac surgery mortality rates: lessons for other specialties. *BMJ*.
237 2013;346:f1139.
- 238 4. Altman DG, Vergouwe Y, Royston P, Moons KGM. Prognosis and prognostic
239 research: validating a prognostic model. *BMJ*. 2009;338(b605):1432-1435.
- 240 5. Harrell Jr FE, Lee KL, Mark DB. Multivariable prognostic models: issues in
241 developing models, evaluating assumptions and adequacy, and measuring
242 and reducing errors. *Stat Med*. 1996;15:361-387.
- 243 6. Steyerberg EW. *Clinical Prediction Models: A Practical Approach to*
244 *Development, Validation, and Updating*. New York: Springer; 2009.
- 245 7. Nilsson J, Ohlsson M, Thulin L, Höglund P, Nashef SAM, Brandt J. Risk factor
246 identification and mortality prediction in cardiac surgery using artificial neural
247 networks. *J Thorac Cardiovasc Surg*. 2006;132(1).
- 248 8. Moons KGM, Altman DG, Vergouwe Y, Royston P. Prognosis and prognostic
249 research: application and impact of prognostic models in clinical practice. *BMJ*.
250 2009;338(b606):1487-1490.
- 251 9. Rufibach K. Use of Brier score to assess binary predictions. *J Clin Epidemiol*.
252 2010;63(8):938-939.
- 253 10. Grant SW, Hickey GL, Dimarakis I, Trivedi U, Bryan AJ, Treasure T, et al. How
254 does EuroSCORE II perform in UK cardiac surgery; an analysis of 23 740
255 patients from the Society for Cardiothoracic Surgery in Great Britain and
256 Ireland National Database. *Heart*. 2012;98(21):1568-1572.
- 257 11. Grant SW, Hickey GL, Dimarakis I, Cooper G, Jenkins DP, Uppal R, et al.

- 258 Performance of the EuroSCORE models in emergency cardiac surgery. *Circ*
259 *Cardiovasc Qual Outcomes*. 2013;6(2):178-185.
- 260 12. Knol MJ, Janssen KJM, Donders RRT, Egberts ACG, Heerdink ER, Grobbee
261 DE, et al. Unpredictable bias when using the missing indicator method or
262 complete case analysis for missing confounder values: an empirical example.
263 *J Clin Epidemiol*. 2010;63(7):728-736.
- 264 13. Marshall A, Altman DG, Holder RL, Royston P. Combining estimates of
265 interest in prognostic modelling studies after multiple imputation: current
266 practice and guidelines. *BMC Med Res Methodol*. 2009;9:57.
- 267 14. Janssen KJM, Vergouwe Y, Donders a. RT, Harrell Jr FE, Chen Q, Grobbee
268 DE, et al. Dealing with missing predictor values when applying clinical
269 prediction models. *Clin Chem*. 2009;55(5):994-1001.
- 270 15. Vergouwe Y, Steyerberg EW, Eijkemans MJC, Habbema JDF. Substantial
271 effective sample sizes were required for external validation studies of
272 predictive logistic regression models. *J Clin Epidemiol*. 2005;58(5):475-483.
- 273 16. Collins GS, Ogundimu EO, Altman DG. Sample size considerations for the
274 external validation of a multivariable prognostic model: a resampling study.
275 *Stat Med*. 2016;35(2):214-226.
- 276 17. Roques F. The logistic EuroSCORE. *Eur Heart J*. 2003;24(9):1-2.
- 277 18. Shahian DM, O'Brien SM, Filardo G, Ferraris VA, Haan CK, Rich JB, et al. The
278 Society of Thoracic Surgeons 2008 cardiac surgery risk models: Part 1 -
279 coronary artery bypass grafting surgery. *Ann Thorac Surg*. 2009;88(1
280 Suppl):S2-S22.
- 281 19. O'Brien SM, Shahian DM, Filardo G, Ferraris VA, Haan CK, Rich JB, et al. The
282 Society of Thoracic Surgeons 2008 cardiac surgery risk models: Part 2 -
283 isolated valve surgery. *Ann Thorac Surg*. 2009;88(1 Suppl):S23-S42.
- 284 20. Shahian DM, O'Brien SM, Filardo G, Ferraris VA, Haan CK, Rich JB, et al. The
285 Society of Thoracic Surgeons 2008 cardiac surgery risk models: Part 3 - valve
286 plus coronary artery bypass grafting surgery. *Ann Thorac Surg*. 2009;88(1
287 Suppl):S43-S62.
- 288 21. Jacobs JP, Mavroudis C, Jacobs ML, Maruszewski B, Tchervenkov CI,

- 289 Lacour-Gayet FG, et al. What is operative mortality? Defining death in a
290 surgical registry database: a report of the STS Congenital Database Taskforce
291 and the Joint EACTS-STC Congenital Database Committee. *Ann Thorac Surg.*
292 2006;81(5):1937-1941.
- 293 22. Swinkels BM, Plokker HW. Evaluating operative mortality of cardiac surgery:
294 first define operative mortality. *Netherlands Hear J.* 2010;18(7-8):344-345.
- 295 23. Nashef SAM, Roques F, Sharples LD, Nilsson J, Smith C, Goldstone AR, et al.
296 EuroSCORE II. *Eur J Cardio-Thoracic Surg.* 2012;41:1-12.
- 297 24. Siregar S, Groenwold RHH, de Mol B a JM, Speekenbrink RGH, Versteegh
298 MIM, Brandon Bravo Bruinsma GJ, et al. Evaluation of cardiac surgery
299 mortality rates: 30-day mortality or longer follow-up? *Eur J Cardio-Thoracic*
300 *Surg.* 2013;44(5):875-883.
- 301 25. Hickey GL, Grant SW, Cosgriff R, Dimarakis I, Pagano D, Kappetein AP, et al.
302 Clinical registries: governance, management, analysis and applications. *Eur J*
303 *Cardio-Thoracic Surg.* 2013;44(4):605-614.
- 304 26. Hickey GL, Grant SW, Murphy GJ, Bhabra M, Pagano D, McAllister K, et al.
305 Dynamic trends in cardiac surgery: why the logistic EuroSCORE is no longer
306 suitable for contemporary cardiac surgery and implications for future risk
307 models. *Eur J Cardio-Thoracic Surg.* 2013;43(6):1146-1152.
- 308 27. Hickey GL, Grant SW, Caiado C, Kendall S, Dunning J, Poullis M, et al.
309 Dynamic prediction modelling approaches for cardiac surgery. *Circ Cardiovasc*
310 *Qual Outcomes.* 2013;6:649-658.
- 311 28. Schulz KF, Altman DG, Moher D, Group C. CONSORT 2010 Statement:
312 updated guidelines for reporting parallel group randomised trials. *Trials.*
313 2010;11(32):1-8.
- 314 29. Moher D, Liberati A, Tetzlaff J, Altman DG, Grp P. Preferred reporting items for
315 systematic reviews and meta-analyses: the PRISMA statement. *PLoS Med.*
316 2009;6(7):1-6.
- 317 30. Collins GS, Reitsma JB, Altman DG, Moons KGM. Transparent reporting of a
318 multivariable prediction model for individual prognosis or diagnosis (TRIPOD):
319 the TRIPOD statement. *Circulation.* 2015;131(2):211-219.

- 320 31. Moons KGM, Altman DG, Reitsma JB, Ioannidis JPA, Macaskill P, Steyerberg
321 EW, et al. Transparent reporting of a multivariable prediction model for
322 individual prognosis Or diagnosis (TRIPOD): explanation and elaboration. *Ann*
323 *Intern Med.* 2015;162(1):W1-W73.
- 324 32. Kramer AA, Zimmerman JE. Assessing the calibration of mortality benchmarks
325 in critical care: The Hosmer-Lemeshow test revisited. *Crit Care Med.*
326 2007;35(9):2052-2056.
- 327 33. Marcin JP, Romano PS. Size matters to a model's fit. *Crit Care Med.*
328 2007;35(9):2212-2213.
- 329 34. Paul P, Pennell L, Lemeshow S. Standardizing the power of the Hosmer-
330 Lemeshow goodness of fit test in large data sets. *Stat Med.* 2013;32:67-80.
- 331 35. Bertolini G, D'Amico R, Nardi D, Tinazzi A, Apolone G. One model, several
332 results: the paradox of the Hosmer-Lemeshow goodness-of-fit test for the
333 logistic regression model. *J Epidemiol Biostat.* 2000;5(4):251-253.
- 334 36. Hosmer DW, Lemeshow S. *Applied Logistic Regression.* Second Edi. New
335 Jersey: John Wiley & Sons, Inc.; 2000.
- 336 37. Austin PC, Steyerberg EW. Graphical assessment of internal and external
337 calibration of logistic regression models by using loess smoothers. *Stat Med.*
338 2014;33(3):517-535.
- 339 38. Harrell Jr FE. *Regression Modeling Strategies: With Applications to Linear*
340 *Models, Logistic Regression, and Survival Analysis.* New York: Springer; 2001.
- 341 39. Cox DR. Two further applications of a model for binary regression. *Biometrika.*
342 1958;45(3-4):562-565.
- 343 40. Spiegelhalter DJ. Probabilistic prediction in patient management and clinical
344 trials. *Stat Med.* 1986;5(5):421-433.
- 345 41. Stallard N. Simple tests for the external validation of mortality prediction
346 scores. *Stat Med.* 2009;28:377-388.
- 347 42. Pencina MJ, D'Agostino Snr RB, D'Agostino Jr RB, Vasan RS. Evaluating the
348 added predictive ability of a new marker: From area under the ROC curve to
349 reclassification and beyond. *Stat Med.* 2008;27:157-172.

- 350 43. Steyerberg EW, Borsboom GJJM, van Houwelingen HC, Eijkemans MJC,
351 Habbema JDF. Validation and updating of predictive logistic regression
352 models: a study on sample size and shrinkage. *Stat Med.* 2004;23(16):2567-
353 2586.
- 354 44. Debray TP, Koffijberg H, Nieboer D, Vergouwe Y, Steyerberg EW, Moons
355 KGM. Meta-analysis and aggregation of multiple published prediction models.
356 *Stat Med.* 2014;33(14):2341-2362.
- 357 45. Su T-L, Jaki T, Hickey GL, Buchan I, Sperrin M. A review of statistical updating
358 methods for clinical prediction models. *Stat Methods Med Res.* 2015:In press.
- 359 46. Harrell Jr FE. rms: Regression Modeling Strategies. 2015. [http://cran.r-](http://cran.r-project.org/package=rms)
360 [project.org/package=rms](http://cran.r-project.org/package=rms).
- 361
- 362

363 **FIGURE LEGEND**

364 **Figure.** A calibration plot for simulated data ($n = 500$). The green triangles denote
365 the mean predicted and observed event probabilities for patients grouped into tenths
366 using deciles. The grey dashed line denotes perfect calibration. A smoothing curve
367 (blue dashed line) and the calibration curve (red solid line) are also overlaid. The
368 distribution of calculated predicted probabilities is overlaid along the horizontal axis.
369 A subset of various statistics useful for validating the model are also shown. This
370 figure was generated using standard statistical software: the rms package for R (R
371 Core Team, R Foundation for Statistical Computing, Vienna, Austria; version 3.1.2).
372 Further details are given in Harrell (2001)³⁸ and Harrell (2015).⁴⁶ Code to reproduce
373 this plot is given in the Appendix.

374

APPENDIX

375

376 **R code to produce figure**

```
377 # If 'rms' package not install, run command
378 # install.packages("rms")
379 library(rms)
380 ## Simulate fake data:
381 ##   y = binary outcome
382 ##   x1, x2, x3 = covariates in the risk model
383 ##   n = sample size
384 set.seed(1)
385 n <- 1000 # 500 development + 500 validation
386 x1 <- runif(n) # covariate 1
387 x2 <- runif(n) # covariate 2
388 x3 <- runif(n) # covariate 3
389 logit <- -5 + 0.5*x1 + 2*x2 + 3.5*x3
390 P <- 1 / (1 + exp(-logit))
391 y <- ifelse(runif(n) <= P, 1, 0) # outcomes
392 d <- data.frame(x1, x2, x3, y) # combined dataset
393
394 ## Fit a risk prediction model to first half of the data
395 f <- lrm(y ~ x1 + x2 + x3, subset = 1:500)
396
397 ## Use model to get predictions for second half of data
398 pred.logit <- predict(f, d[501:1000, ])
399 phat <- 1 / (1 + exp(-pred.logit))
400
401 ## Validate prediction
402 val.prob(phat, y[501:1000], g = 10, riskdist = "predicted")
```