

Supplementary Material: Model Selection in Semiparametric Expectile Regression

Elmar Spiegel^{a*}, Fabian Sobotka^b, Thomas Kneib^a

^aGeorg-August-University Göttingen, Germany

^bCarl von Ossietzky University Oldenburg, Germany

This is the supplementary material of the paper *Approaches for Model Selection in Semiparametric Expectile Regression*.

In the following first the results of the simulation study are given as pictures. Afterwards the selected models of the application are given as tables. The description of the simulation study is given in the paper in Section 4.1. For the design of the application have look at Section 5 of the paper.

The results of the simulation study (Part I) are structured as follows: (1) The results with 2000 observation are given, (2) the results with 500 observations are given. Inside these blocks the different selection approaches concerning P-splines come in the ordering: (i) Decomposition into linear and nonlinear effect ("complete"), (ii) Nonlinear vs. linear vs. no effect ("restricted"), (iii) Nonlinear vs. no effect ("no"), (iv) linear vs. no effect ("parametric") (see Section 3.1 of the paper for the details of this selection approaches). For all these selection approaches the three different data designs ((a) parallel, (b) linear, (c) exponential) are available.

The best model for the nutritional status (Part II) is selected via (i) cross-validation and scoring, (ii) stepwise AIC and area under the AIC curve, (iii) non-negative garrote and non-negative garrote on the grid and (iv) Boosting. Besides the results of weighted scoring are given in Table A.7. Furthermore the results of the selection of relevant effects per asymmetry parameter via stepwise backwards CV selection are given in Table A.8. Finally the estimated regional effects of the model selected via scoring are plotted in Figure A.25.

Contents

I	Simulation study	3
1	n=2000	3
	i.) Decomposition into linear and nonlinear effect	3
	(a) Parallel design	3
	(b) Linear design	5
	(c) Exponential design	6
	ii.) Nonlinear vs. linear vs. no effect	7
	(a) Parallel design	7
	(b) Linear design	8
	(c) Exponential design	9
	iii.) Nonlinear vs. no effect	10

*Corresponding author: Elmar Spiegel (espiege@uni-goettingen.de), Chair of Statistics, Georg-August-University Göttingen, Humboldtallee 3, 37073 Göttingen, Germany

(a)	Parallel design	10
(b)	Linear design	11
(c)	Exponential design	12
iv.)	Linear vs. no effect	13
(a)	Parallel design	13
(b)	Linear design	14
(c)	Exponential design	15
2	n=500	16
i.)	Decomposition into linear and nonlinear effect	16
(a)	Parallel design	16
(b)	Linear design	18
(c)	Exponential design	19
ii.)	Nonlinear vs. linear vs. no effect	20
(a)	Parallel design	20
(b)	Linear design	21
(c)	Exponential design	22
iii.)	Nonlinear vs. no effect	23
(a)	Parallel design	23
(b)	Linear design	24
(c)	Exponential design	25
iv.)	Linear vs. no effect	26
(a)	Parallel design	26
(b)	Linear design	27
(c)	Exponential design	28
II	Application	29
3	Selection results for all approaches	29
i.)	Cross-validation and scoring	29
ii.)	Stepwise AIC and area under the AIC curve	30
iii.)	Non-negative garrote and non-negative garrote on the grid	31
iv.)	Expectile Boosting	32
v.)	Quantile Boosting	33
vi.)	Selection based on confidence intervals	34
vii.)	Standard scoring vs. weighted scoring	35
viii.)	Stepwise backward CV after scoring	36
4	Estimated regional effects for all used asymmetries	37

Part I

Simulation study

1 $n=2000$

i.) Decomposition into linear and nonlinear effect

(a) Parallel design

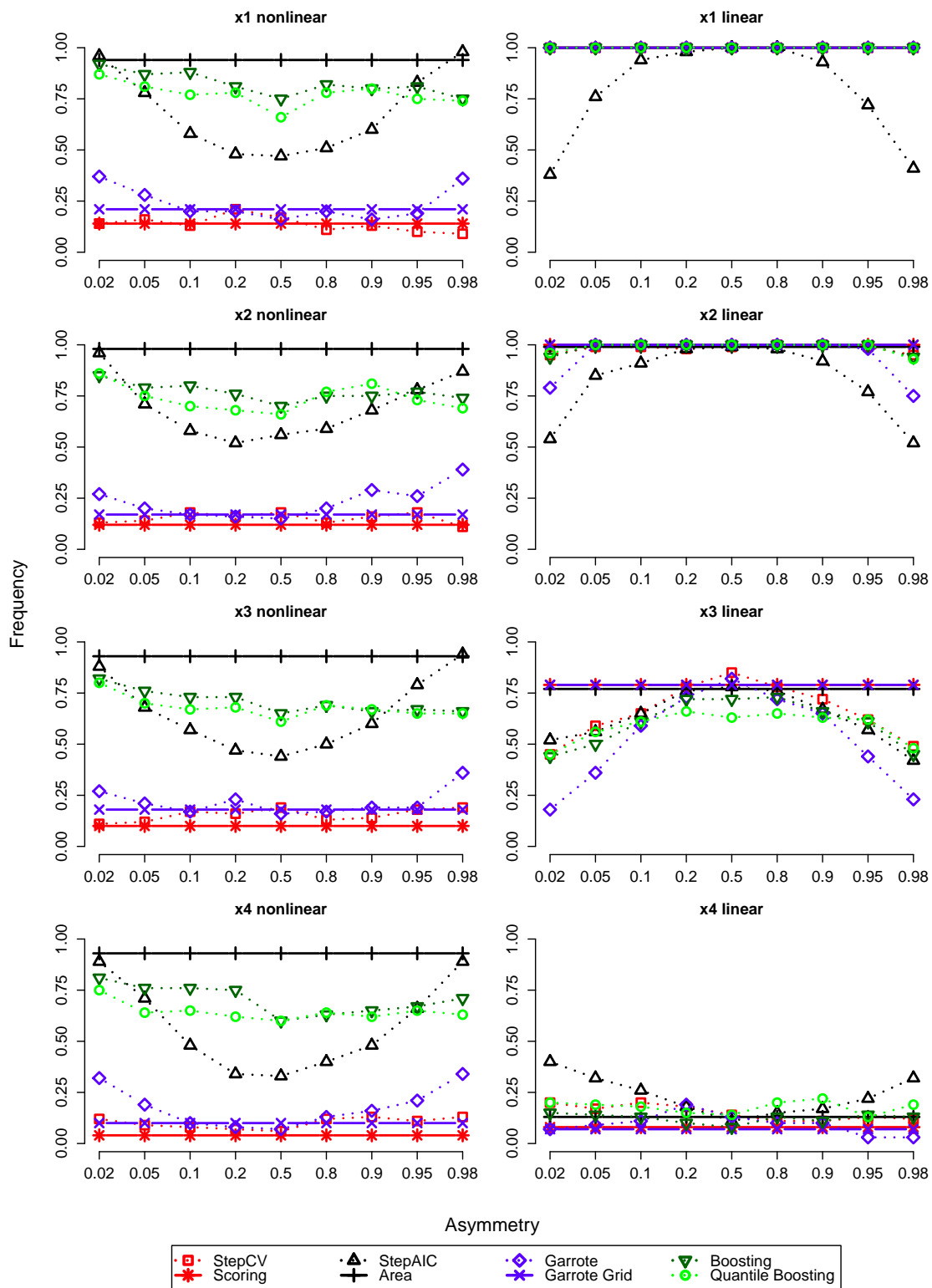


Figure A.1: Frequency of selected models for parallel design with $n=2000$ and decomposition into linear and nonlinear effect

(b) Linear design

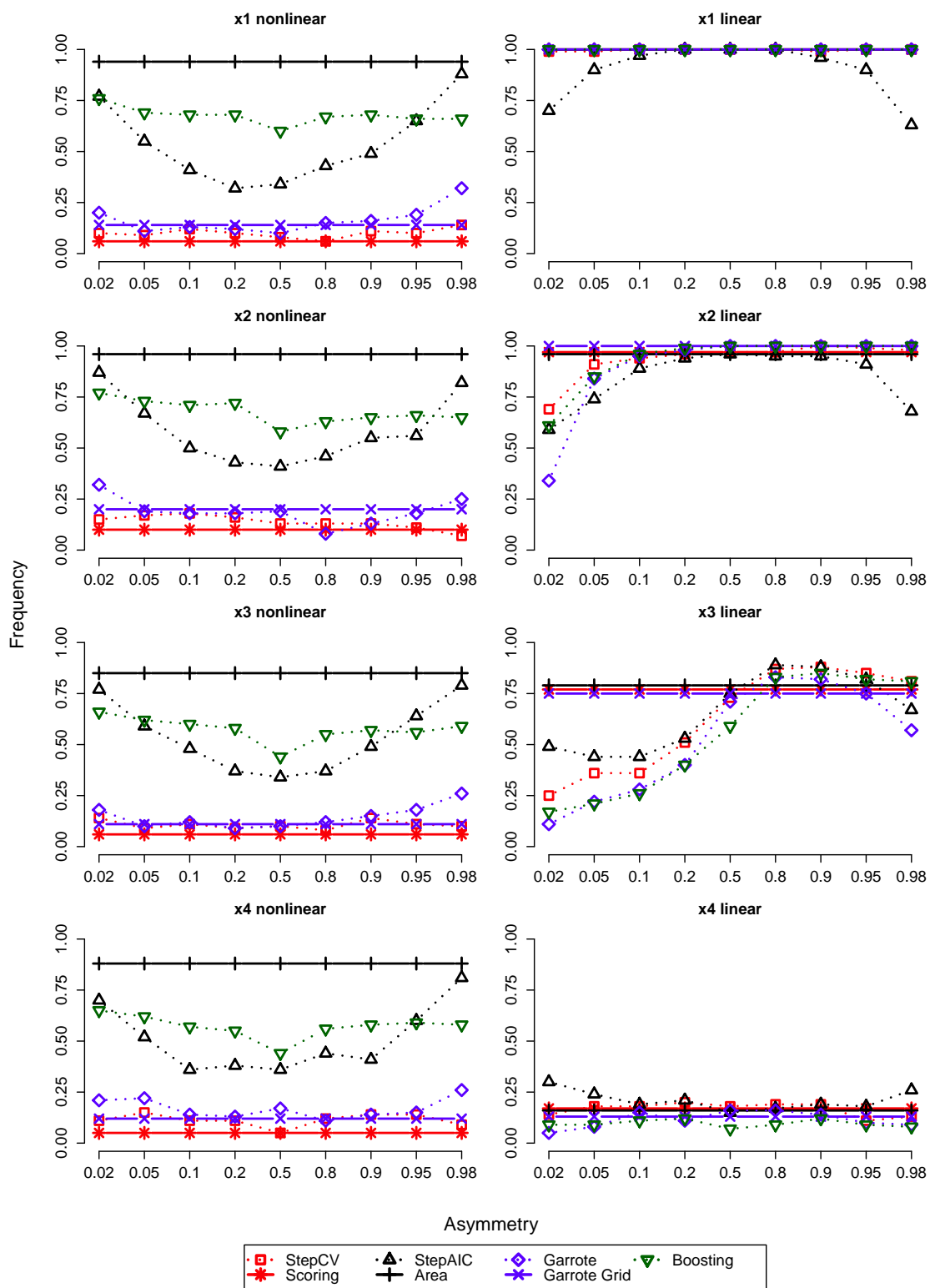


Figure A.2: Frequency of selected models for linear design with $n=2000$ and decomposition into linear and nonlinear effect

(c) Exponential design

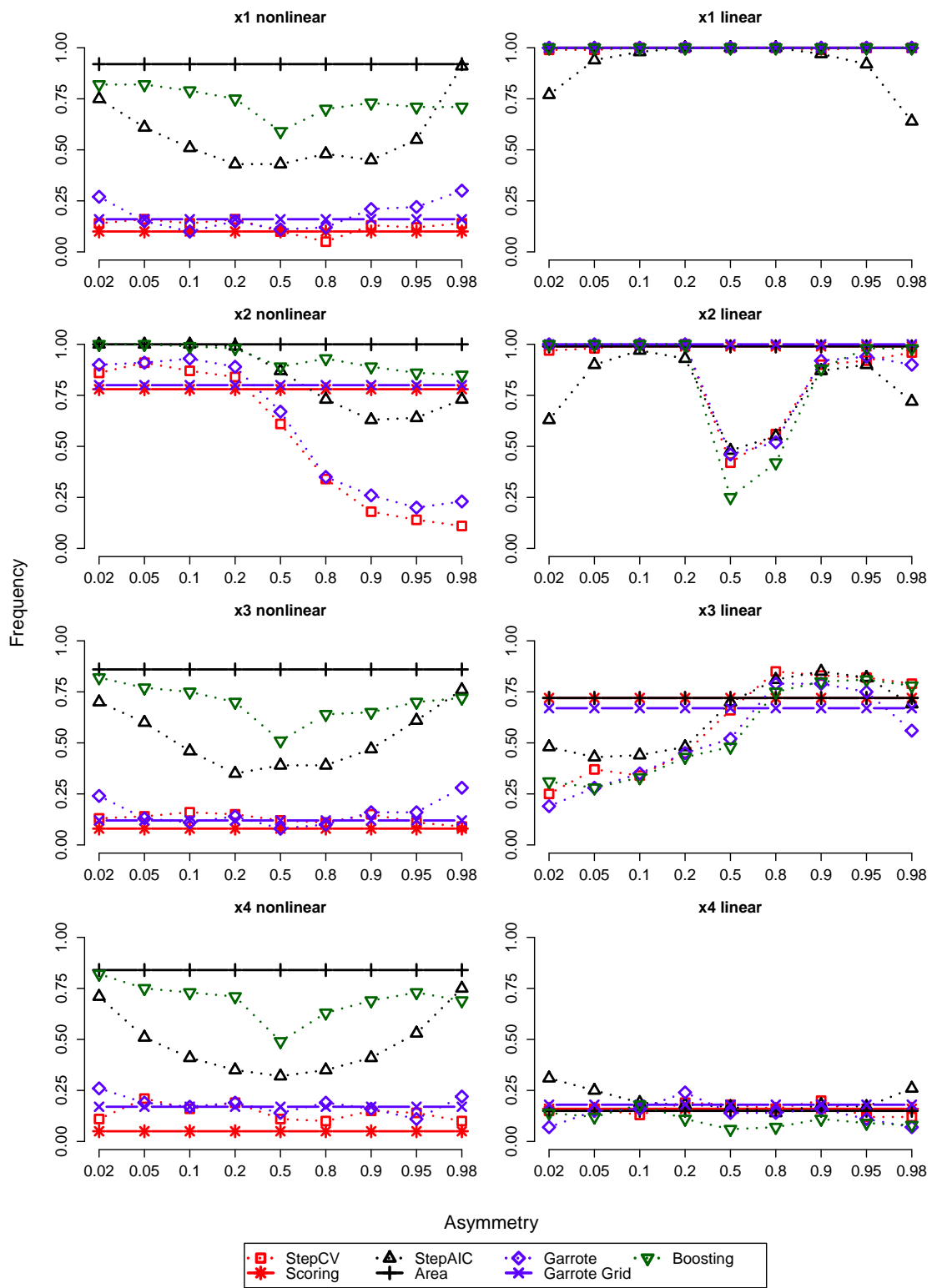


Figure A.3: Frequency of selected models for exponential design with $n=2000$ and decomposition into linear and nonlinear effect

ii.) Nonlinear vs. linear vs. no effect

(a) Parallel design

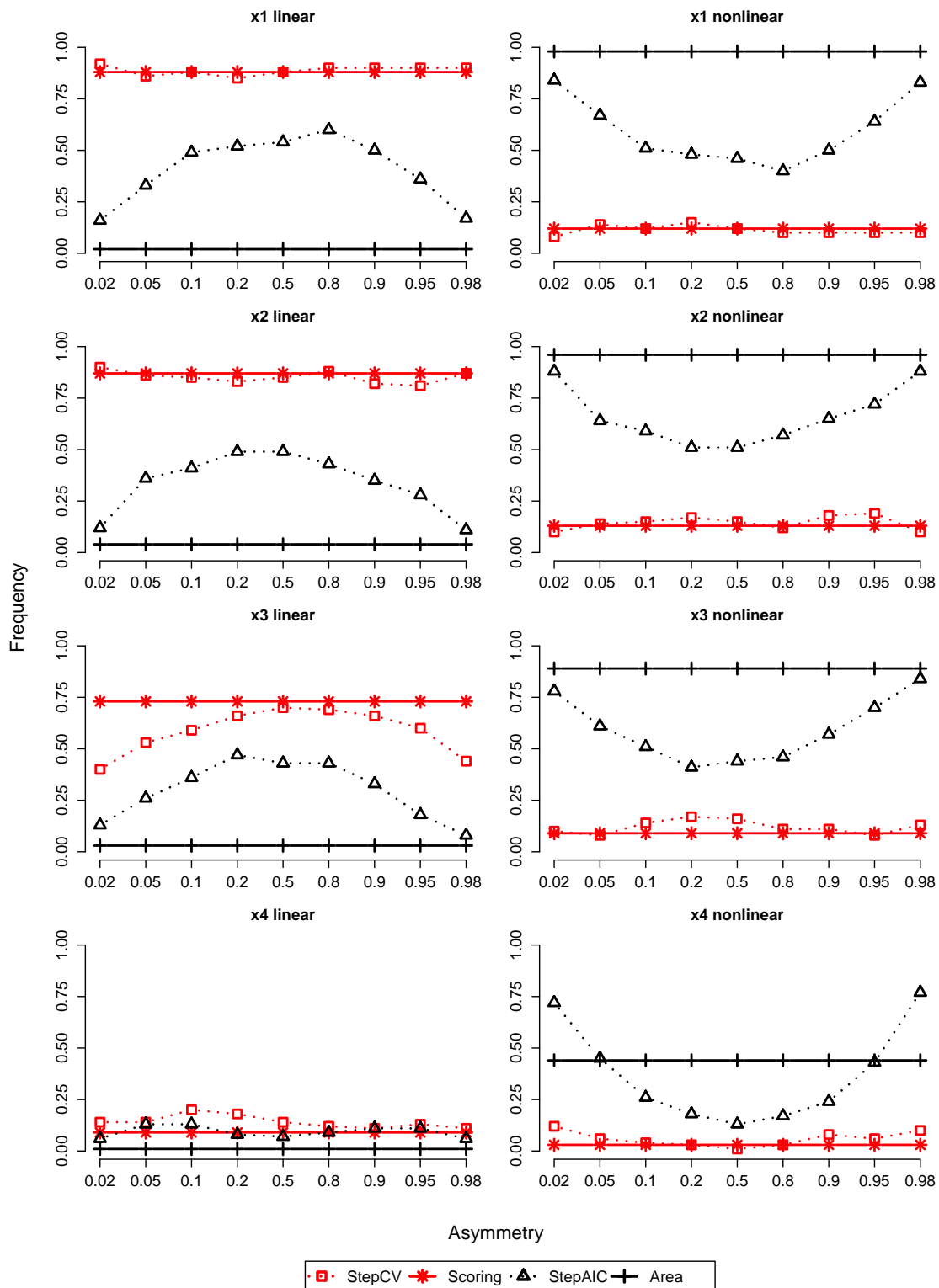


Figure A.4: Frequency of selected models for parallel design with $n=2000$ and restricted selection of P-splines

(b) Linear design

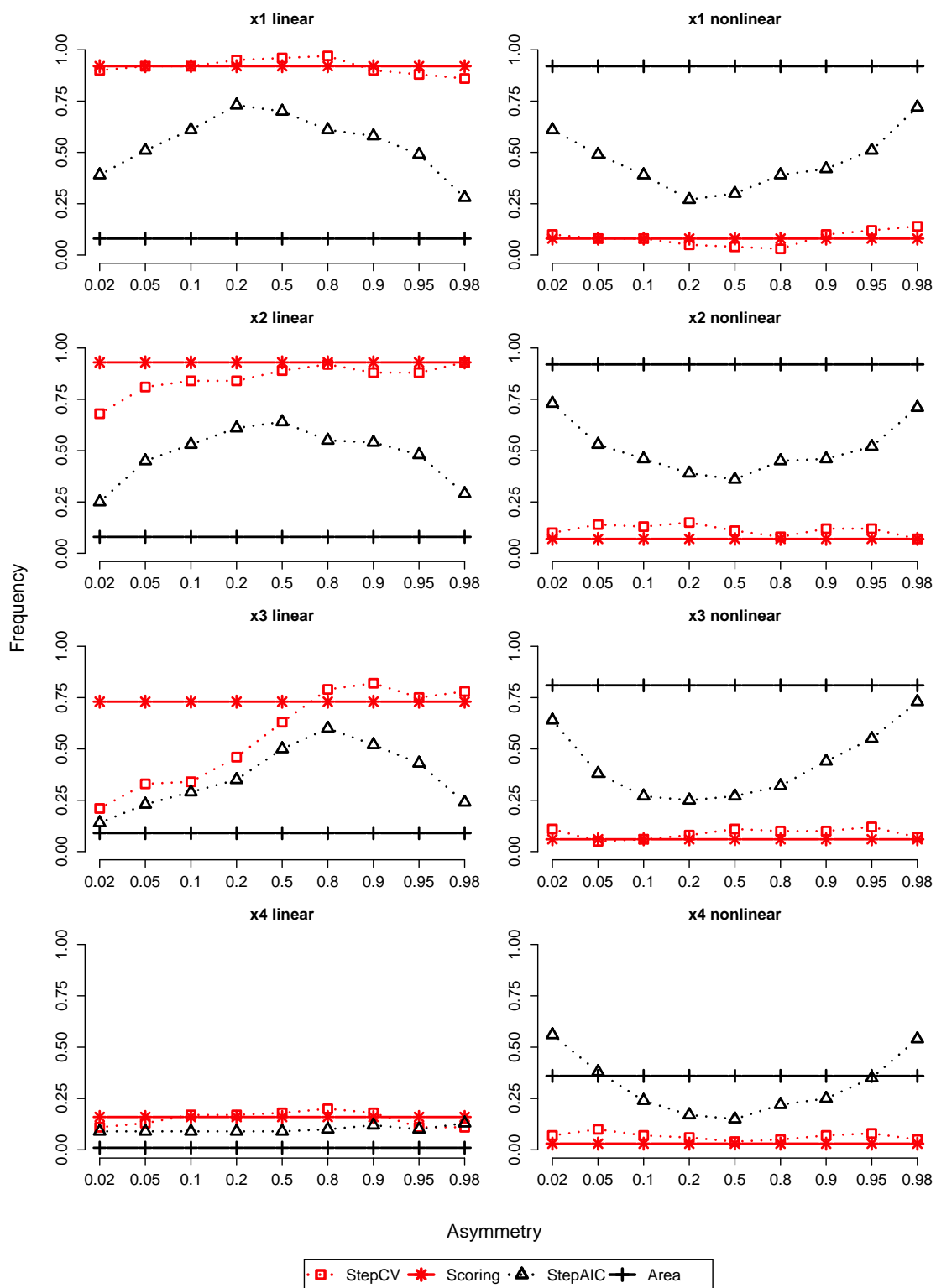


Figure A.5: Frequency of selected models for linear design with $n=2000$ and restricted selection of P-splines

(c) Exponential design

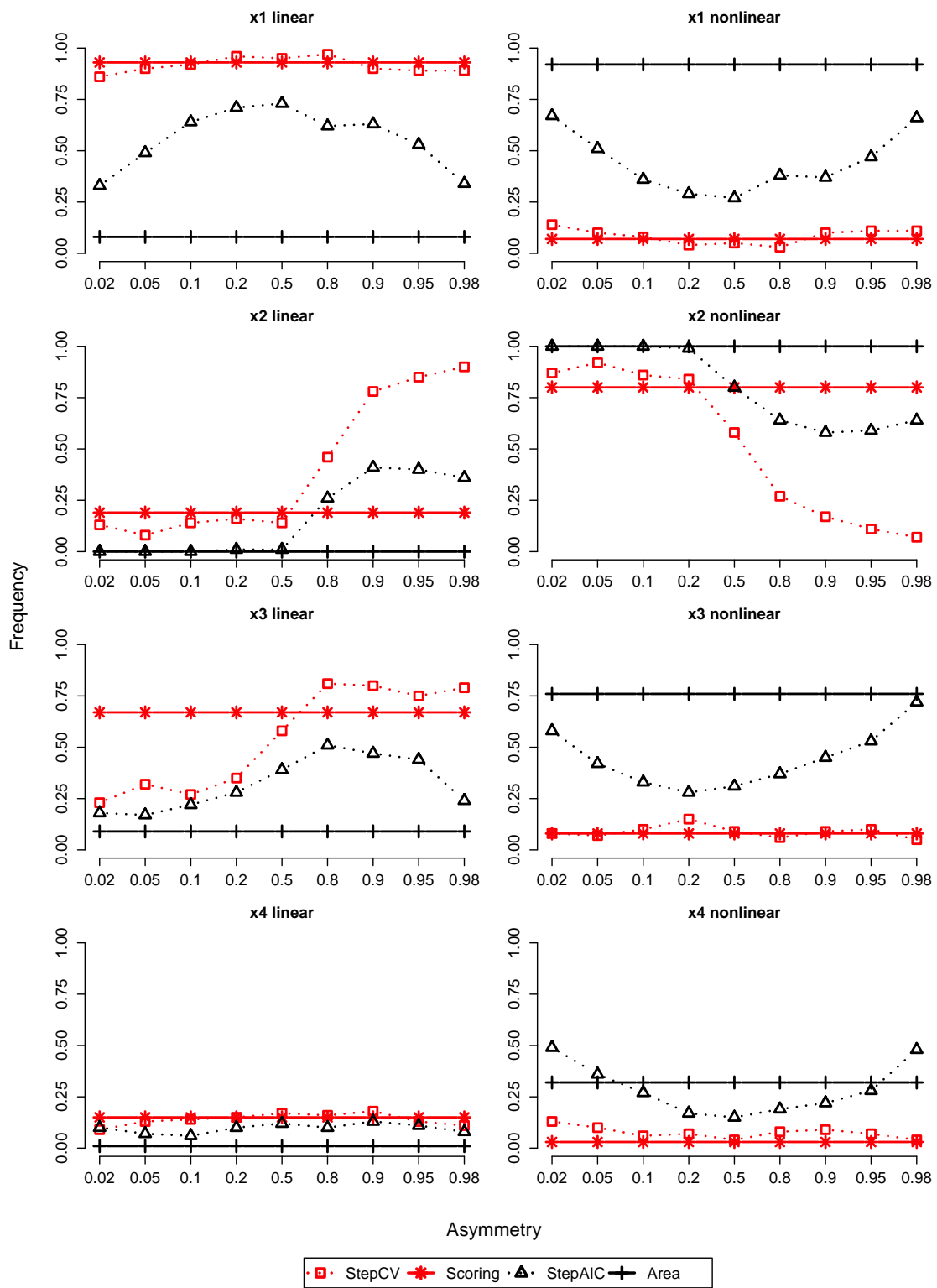


Figure A.6: Frequency of selected models for exponential design with $n=2000$ and restricted selection of P-splines

iii.) Nonlinear vs. no effect

(a) Parallel design

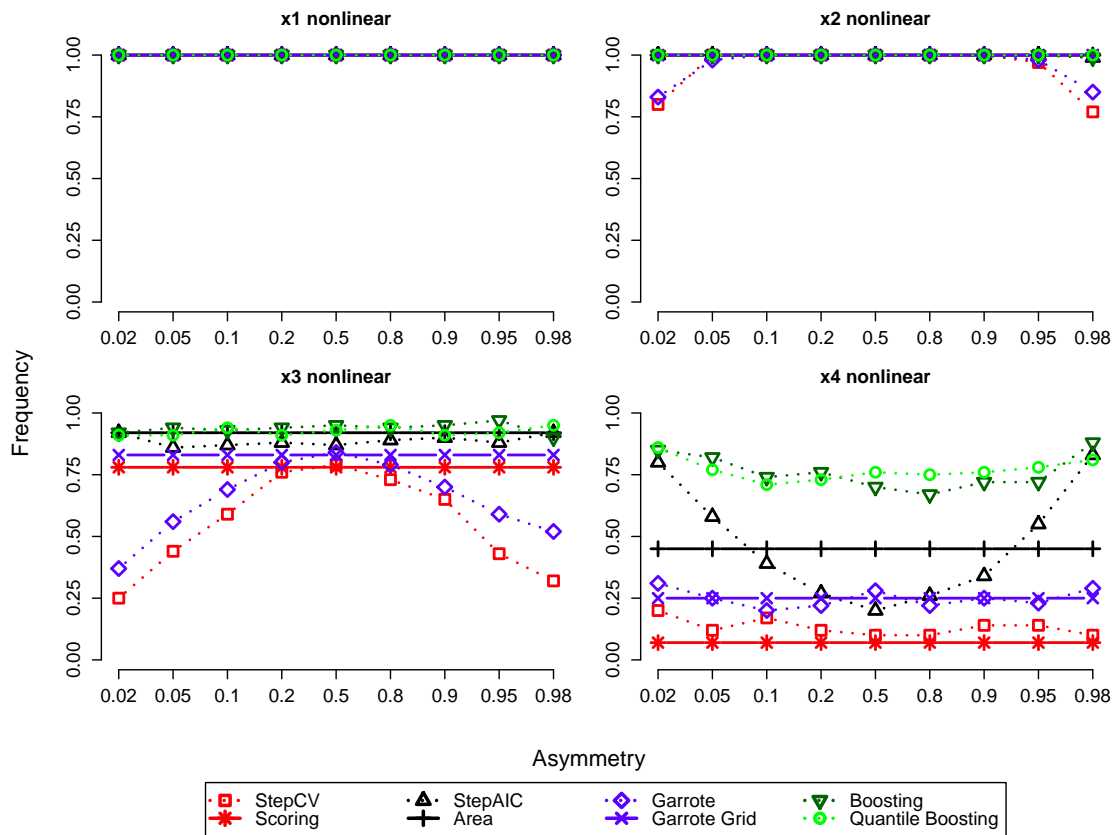


Figure A.7: Frequency of selected models for parallel design with $n=2000$ and selection as nonlinear or no effect

(b) Linear design

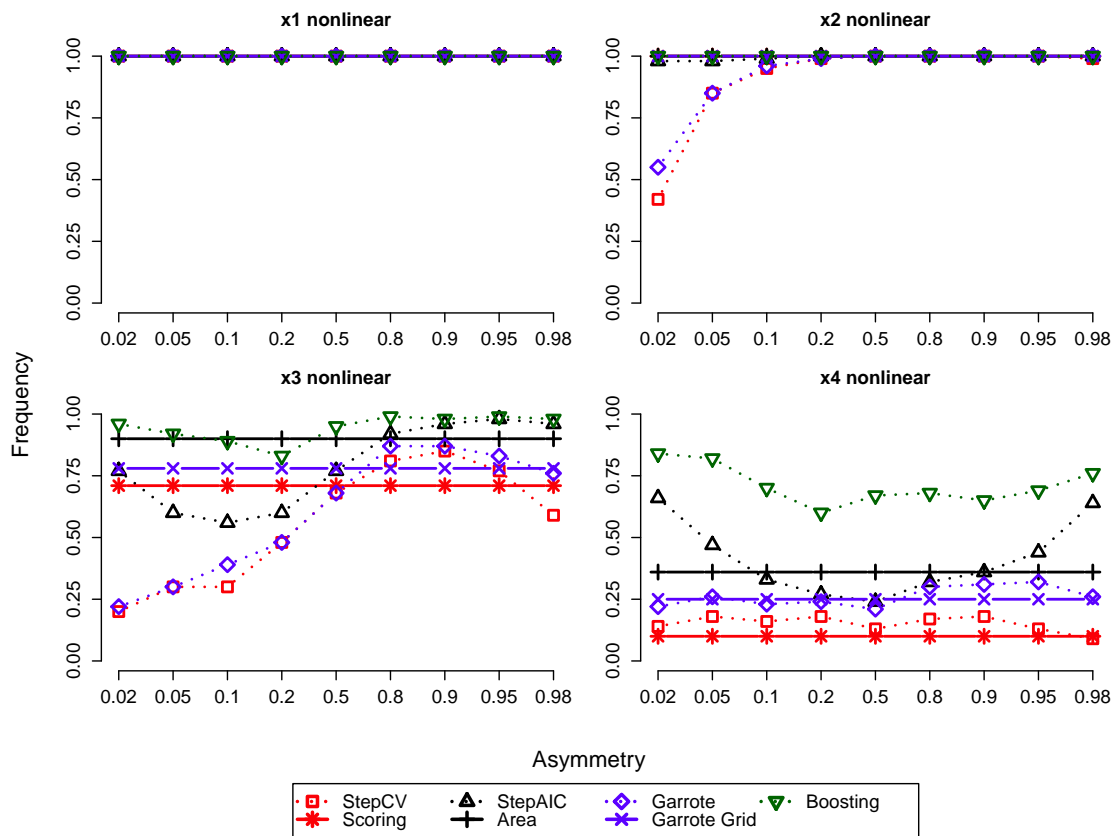


Figure A.8: Frequency of selected models for linear design with $n=2000$ and selection as nonlinear or no effect

(c) Exponential design

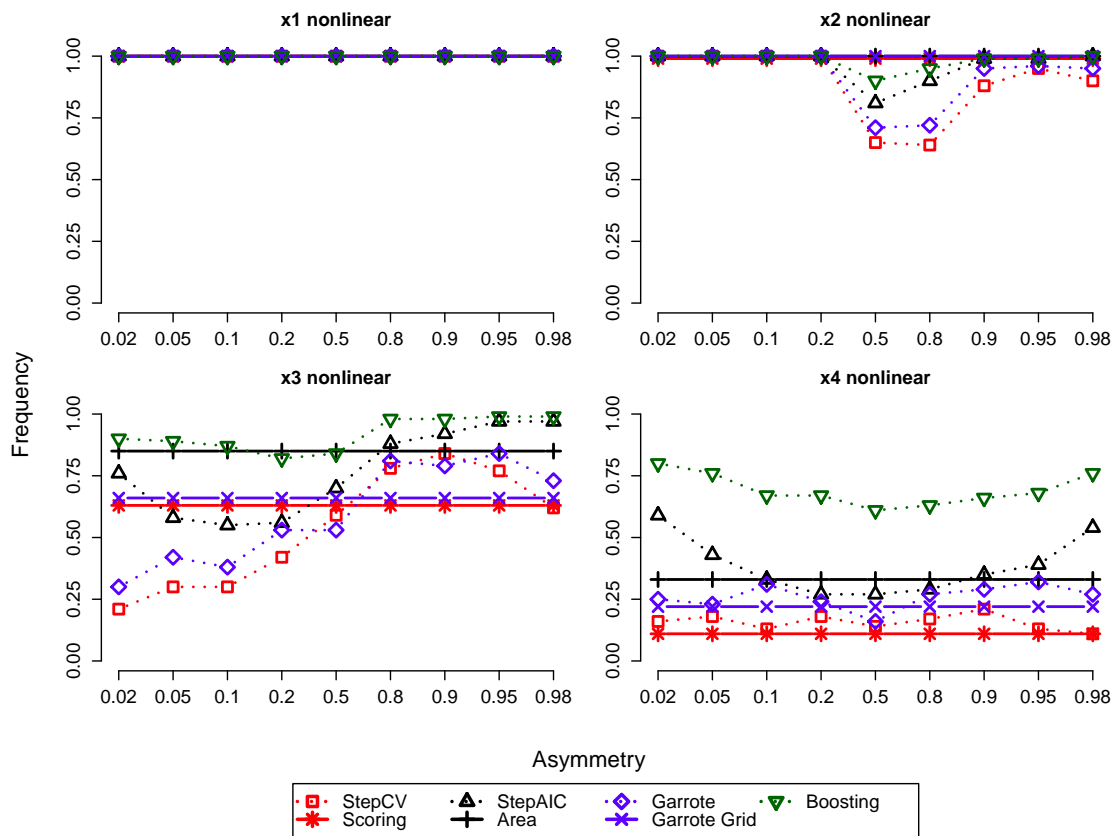


Figure A.9: Frequency of selected models for exponential design with $n=2000$ and selection as nonlinear or no effect

iv.) Linear vs. no effect

(a) Parallel design

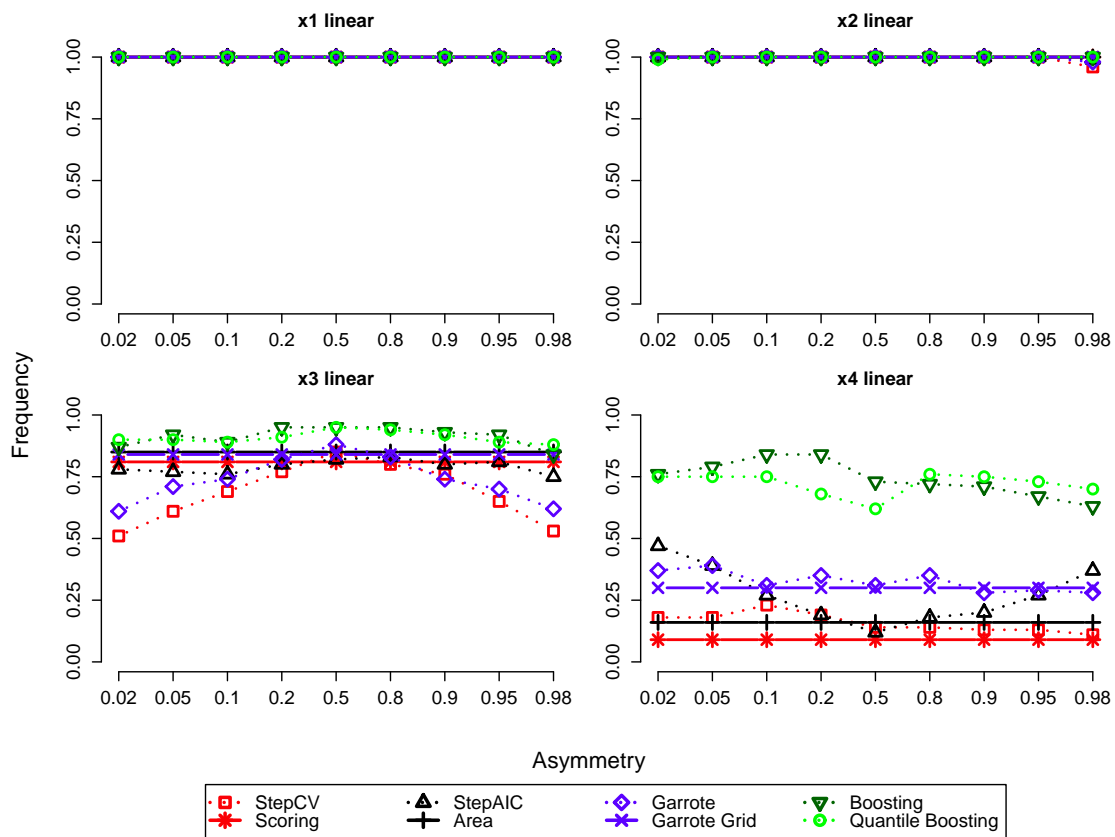


Figure A.10: Frequency of selected models for parallel design with $n=2000$ and selection as linear or no effect

(b) Linear design

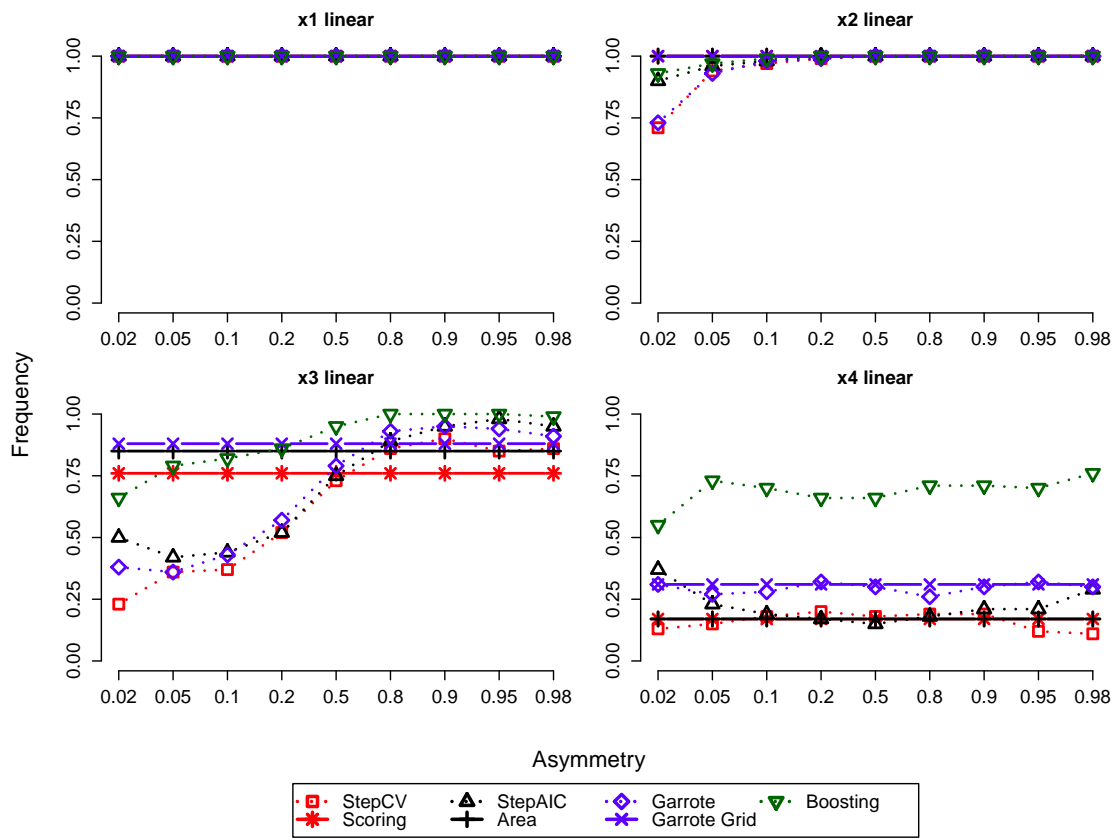


Figure A.11: Frequency of selected models for linear design with $n=2000$ and selection as linear or no effect

(c) Exponential design

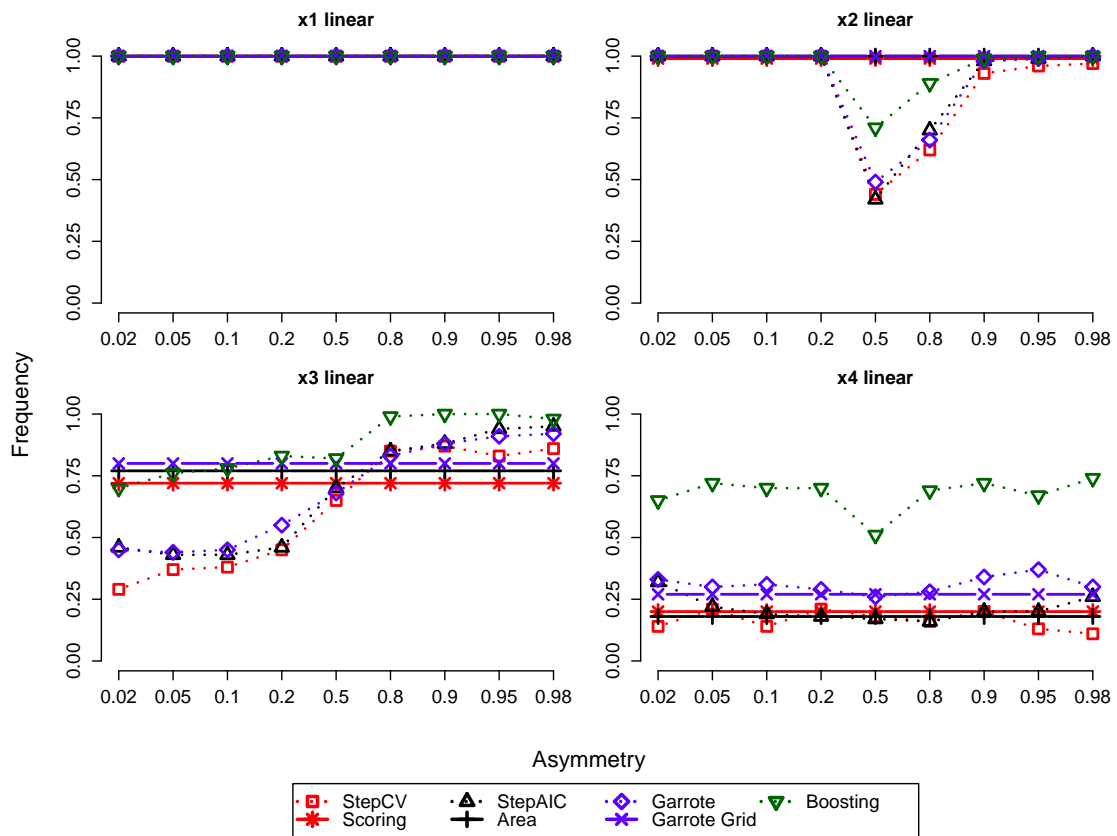


Figure A.12: Frequency of selected models for exponential design with $n=2000$ and selection as linear or no effect

2 n=500

i.) Decomposition into linear and nonlinear effect

(a) Parallel design

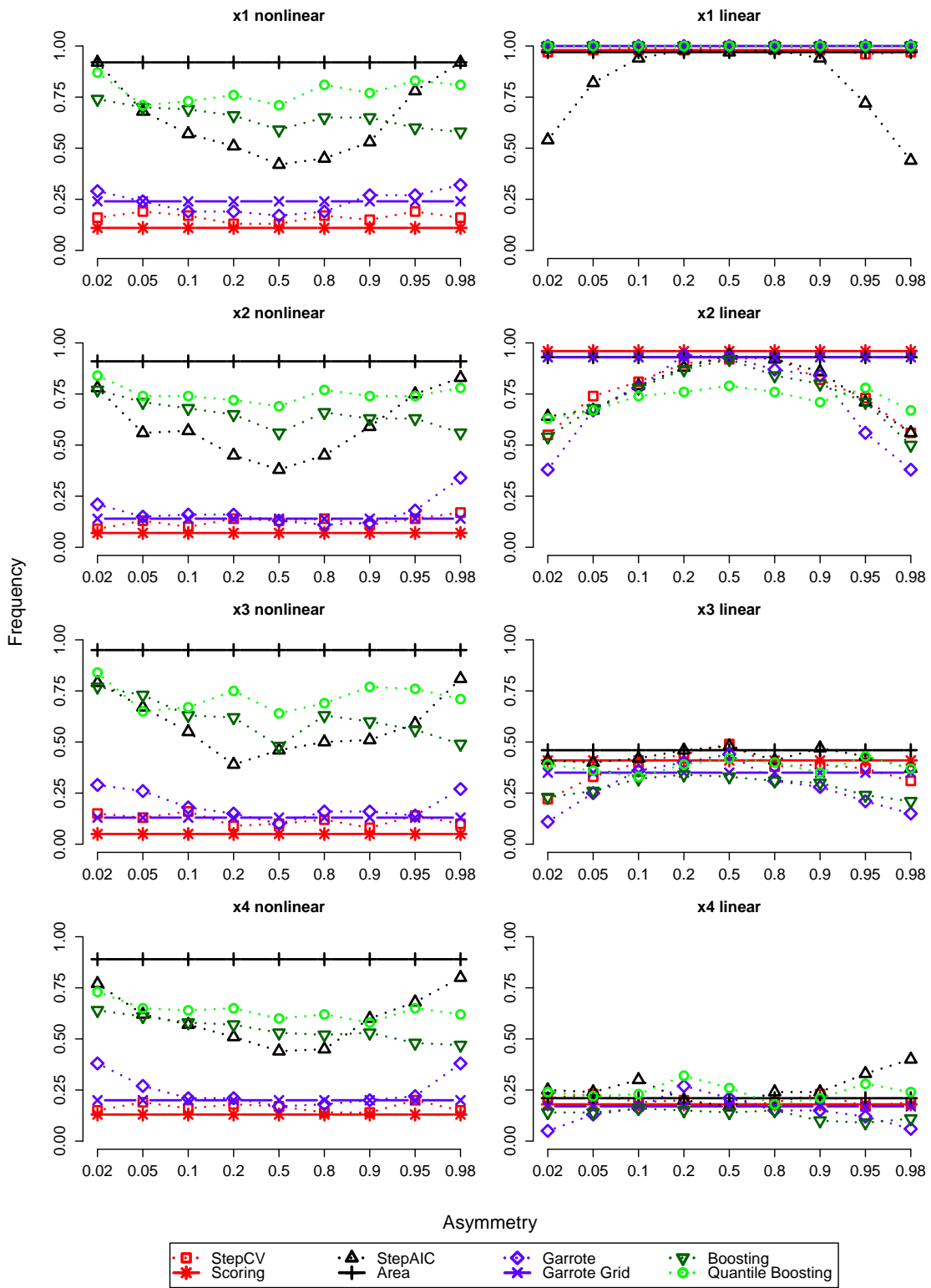


Figure A.13: Frequency of selected models for parallel design with $n=500$ and decomposition into linear and nonlinear effect

(b) Linear design

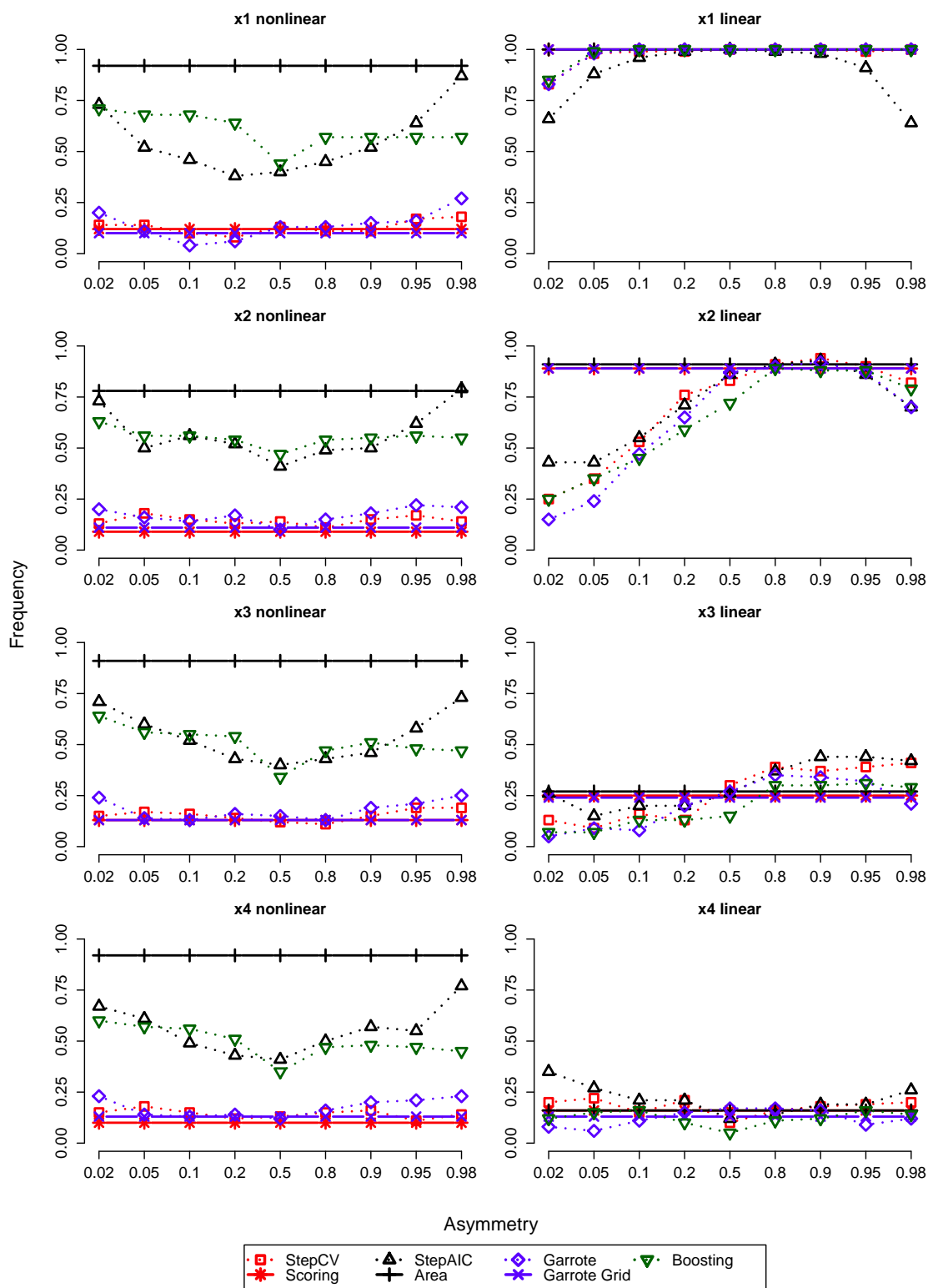


Figure A.14: Frequency of selected models for linear design with $n=500$ and decomposition into linear and nonlinear effect

(c) Exponential design

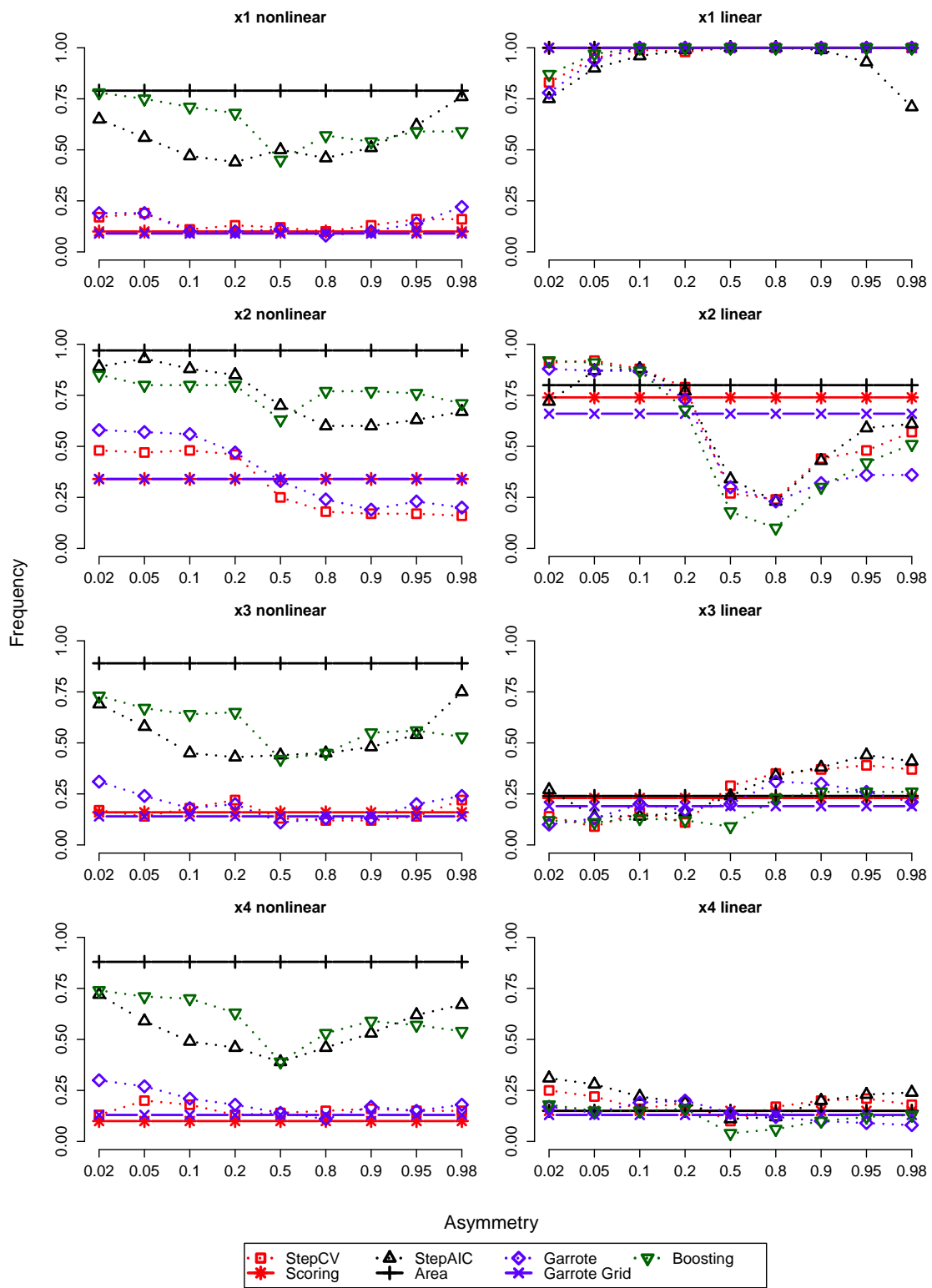


Figure A.15: Frequency of selected models for exponential design with $n=500$ and decomposition into linear and nonlinear effect

ii.) Nonlinear vs. linear vs. no effect

(a) Parallel design

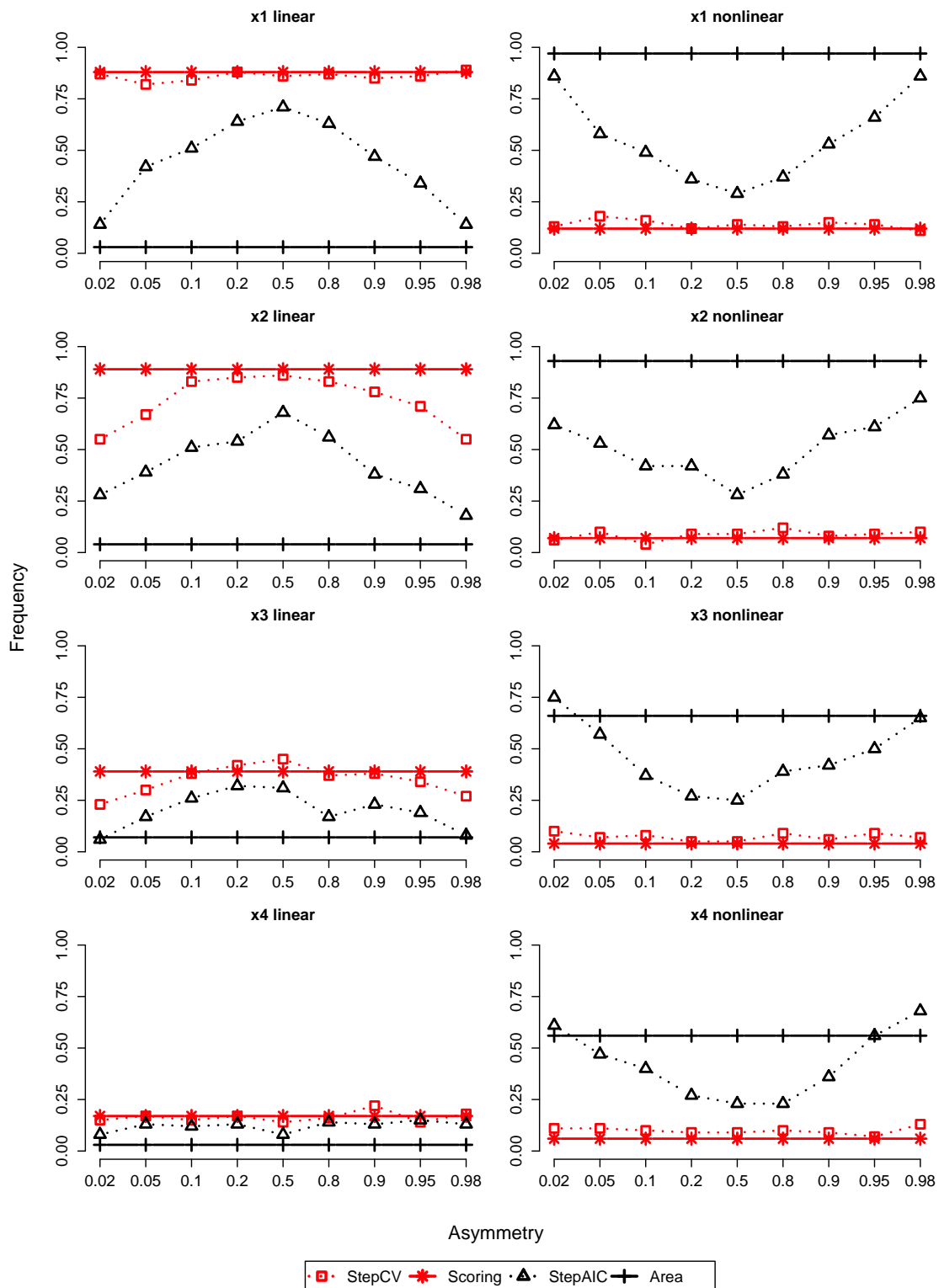


Figure A.16: Frequency of selected models for parallel design with $n=500$ and restricted selection of P-splines

(b) Linear design

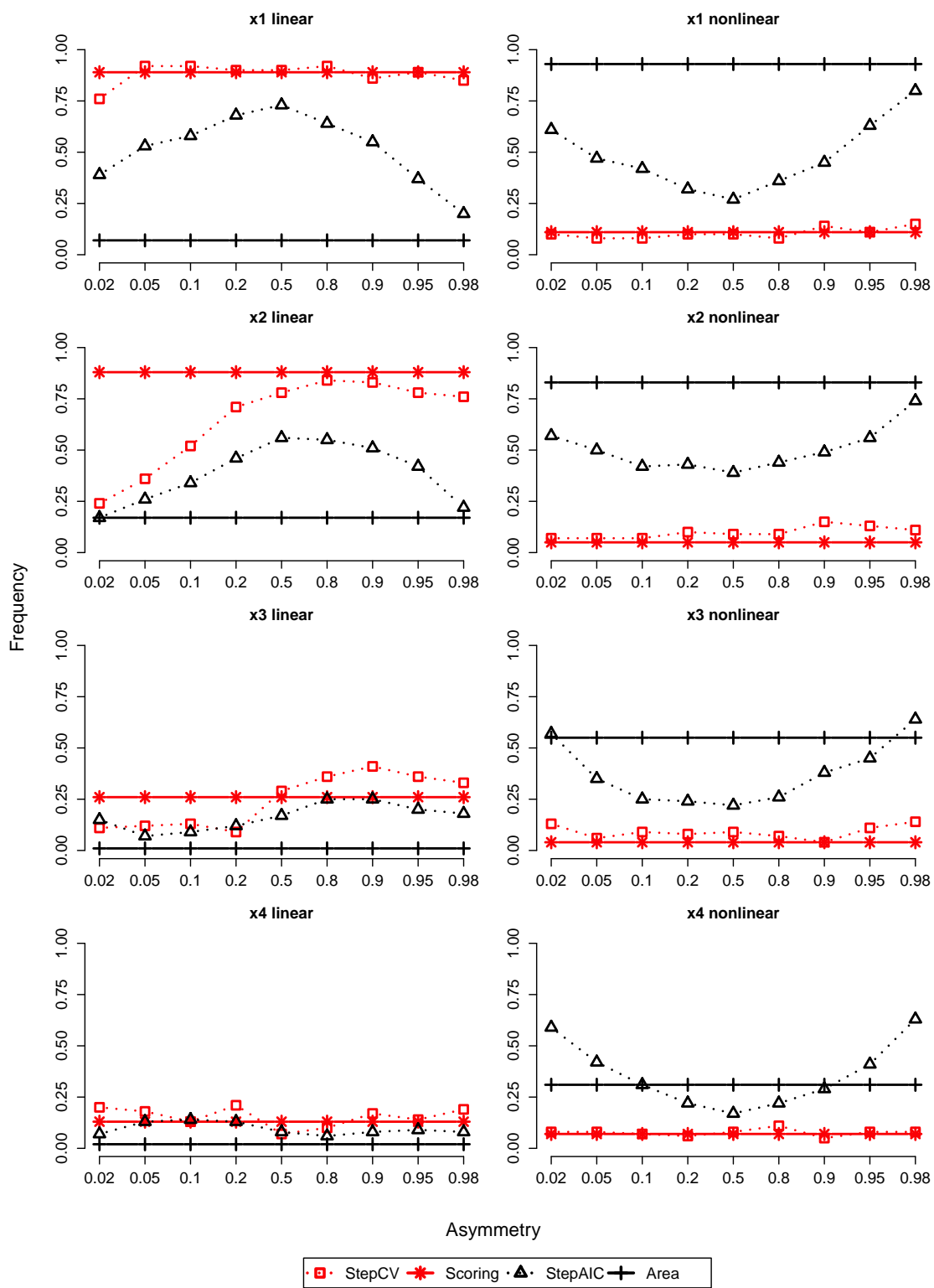


Figure A.17: Frequency of selected models for linear design with $n=500$ and restricted selection of P-splines

(c) Exponential design

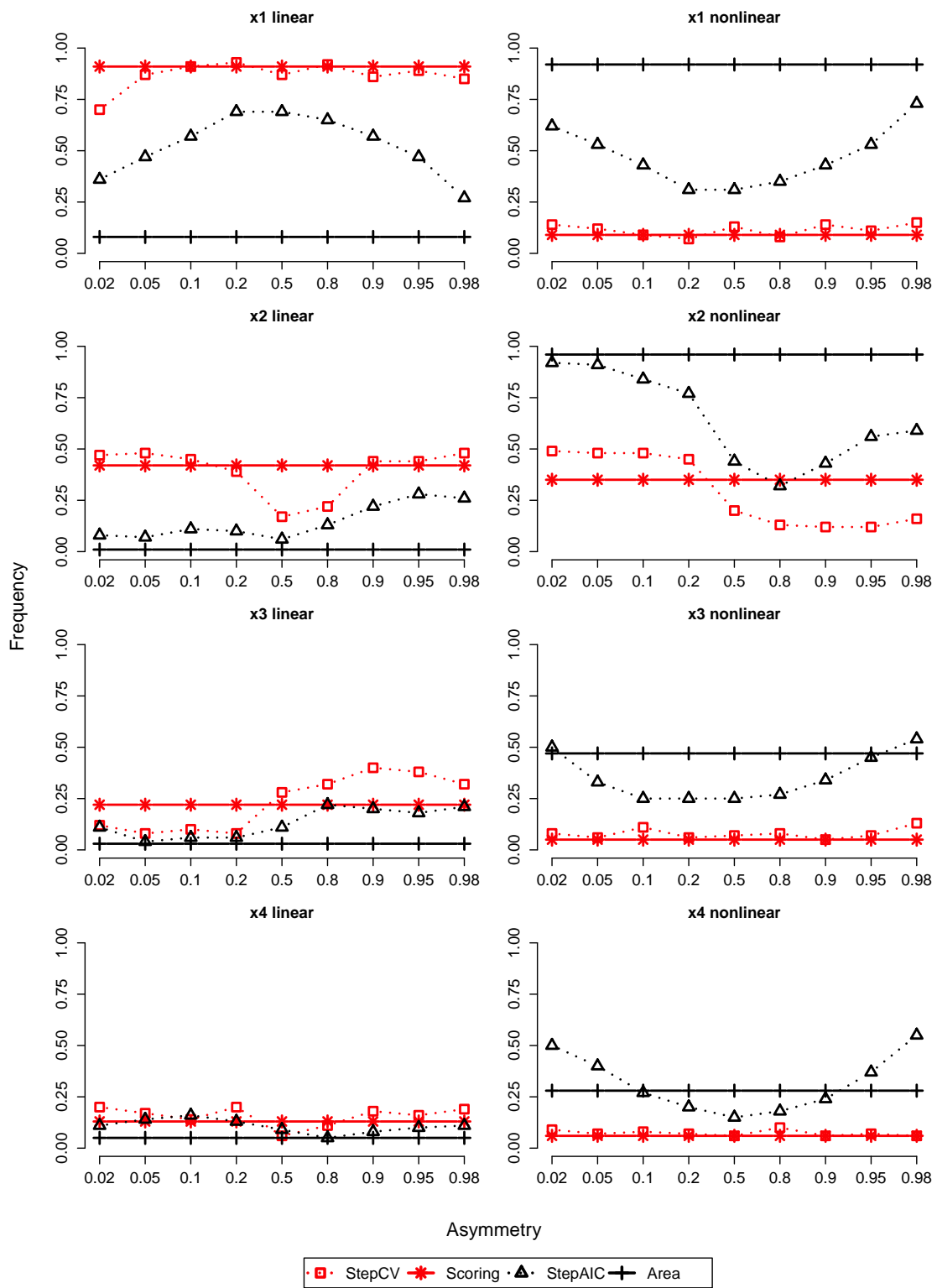


Figure A.18: Frequency of selected models for exponential design with $n=500$ and restricted selection of P-splines

iii.) Nonlinear vs. no effect

(a) Parallel design

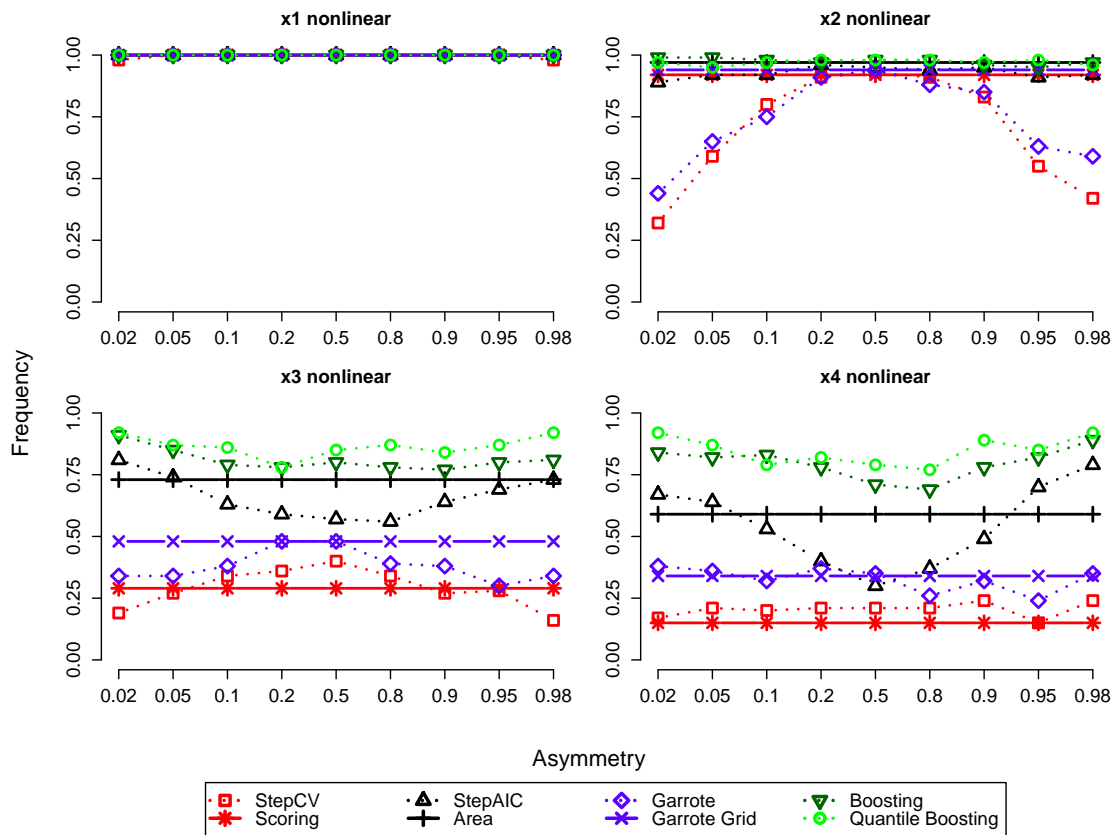


Figure A.19: Frequency of selected models for parallel design with $n=500$ and selection as nonlinear or no effect

(b) Linear design

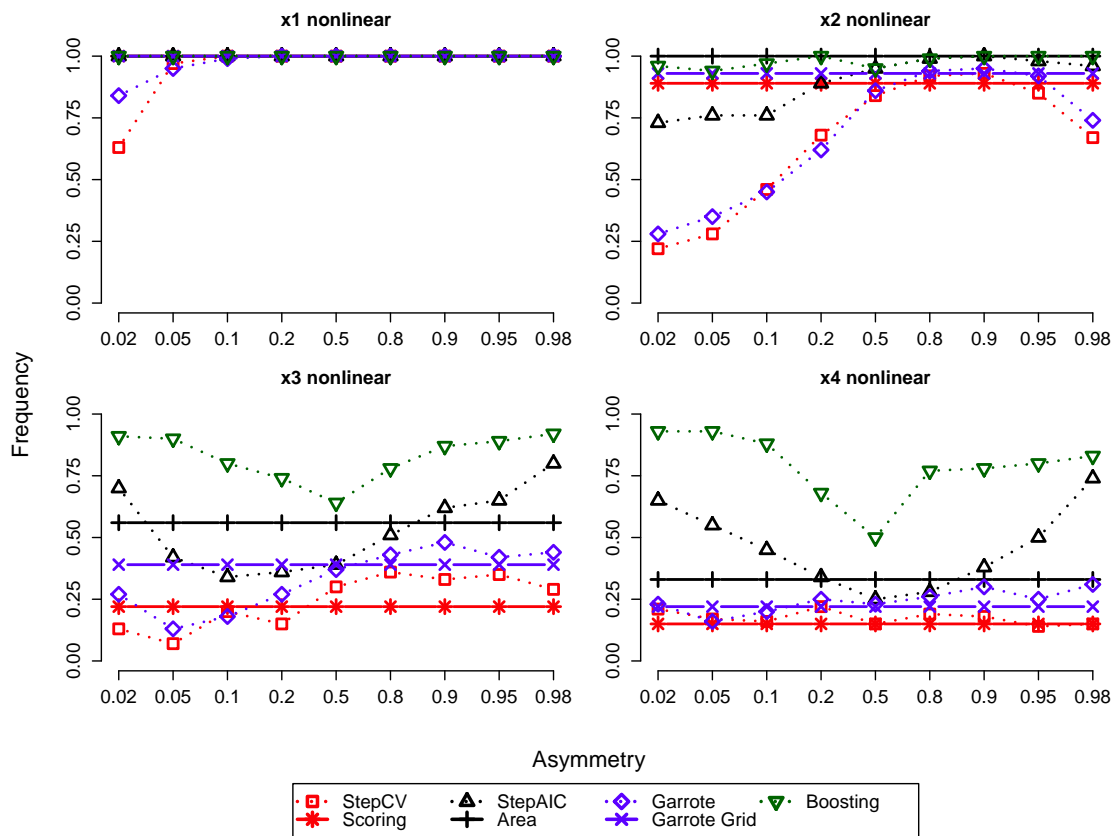


Figure A.20: Frequency of selected models for linear design with $n=500$ and selection as nonlinear or no effect

(c) Exponential design

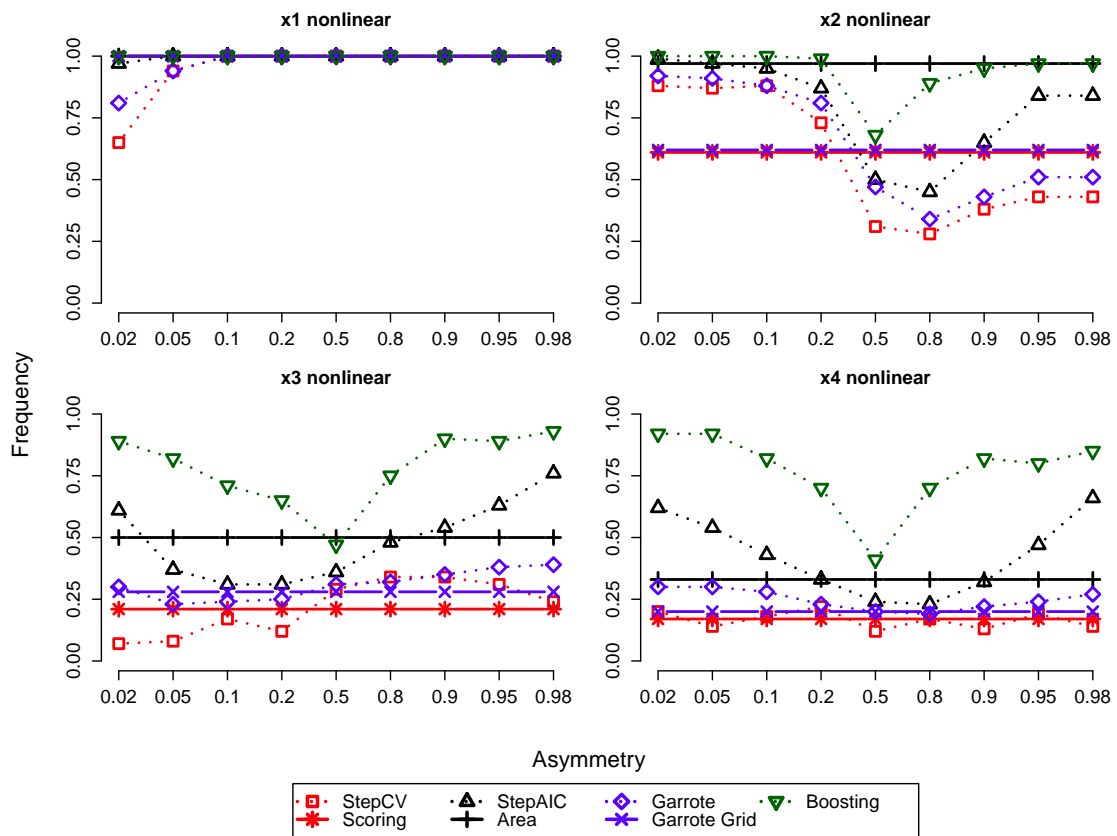


Figure A.21: Frequency of selected models for exponential design with $n=500$ and selection as nonlinear or no effect

iv.) Linear vs. no effect

(a) Parallel design

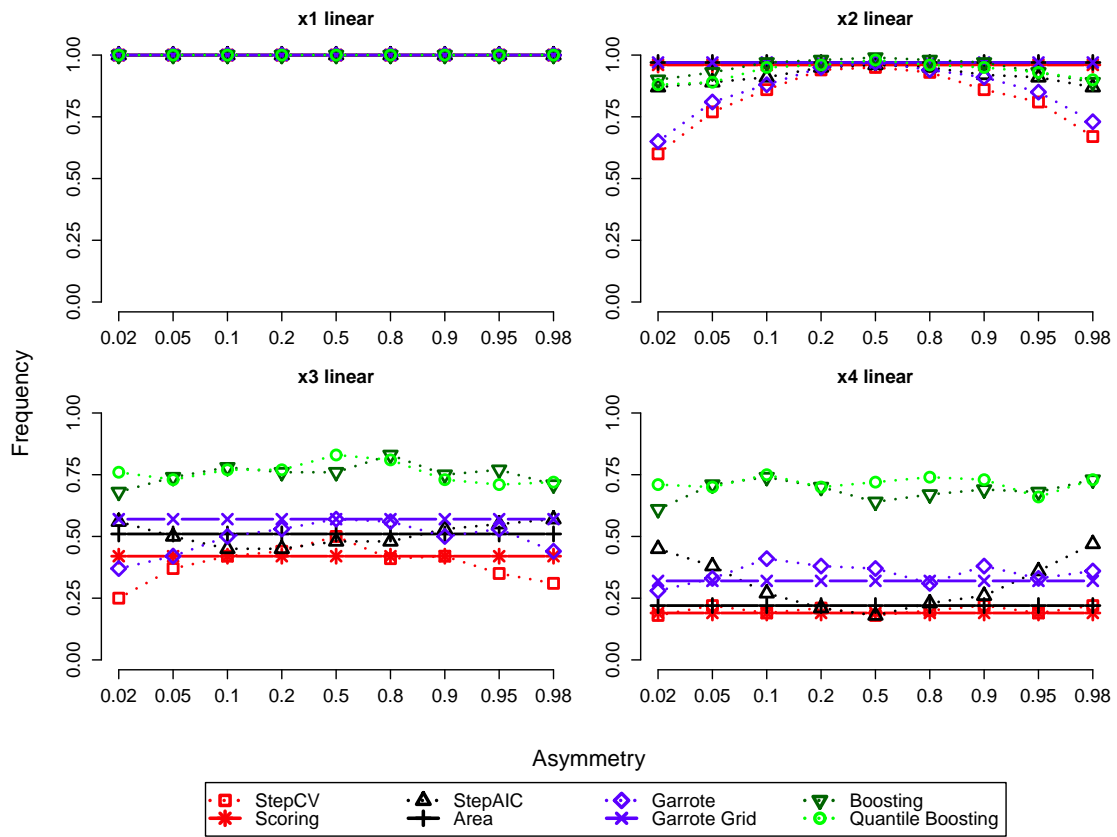


Figure A.22: Frequency of selected models for parallel design with n=500 and selection as linear or no effect

(b) Linear design

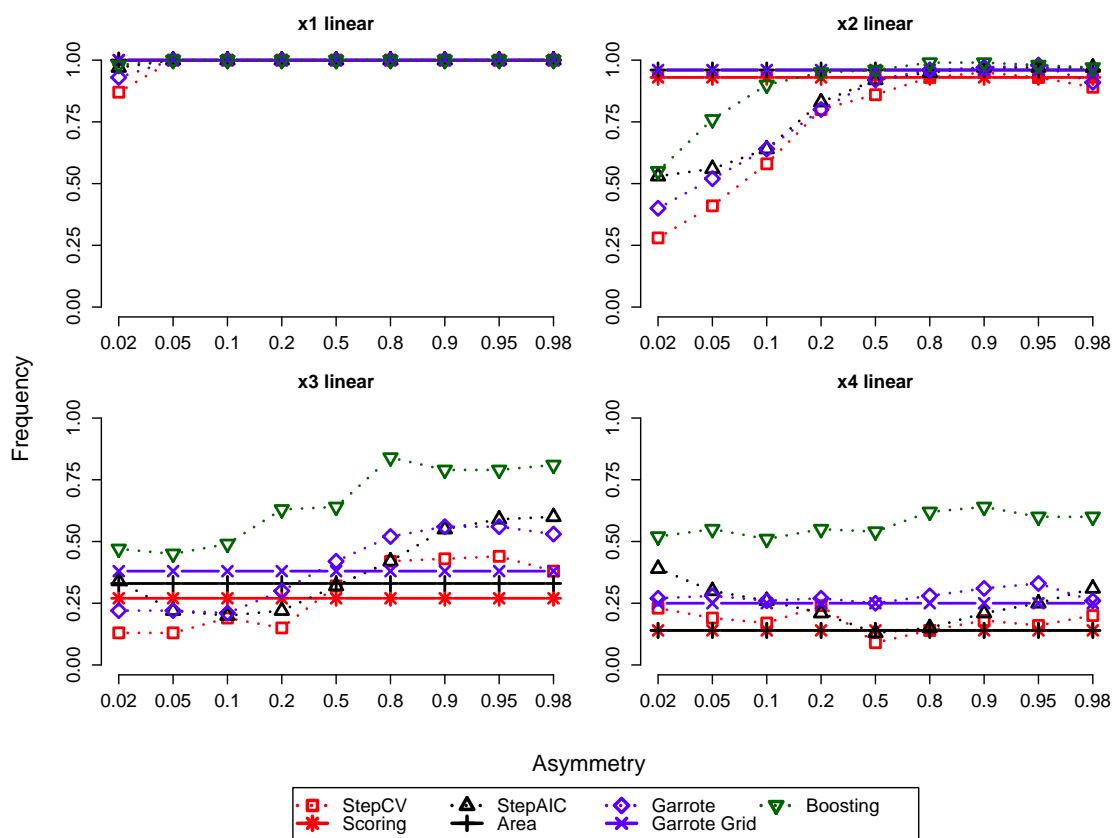


Figure A.23: Frequency of selected models for linear design with $n=500$ and selection as linear or no effect

(c) Exponential design

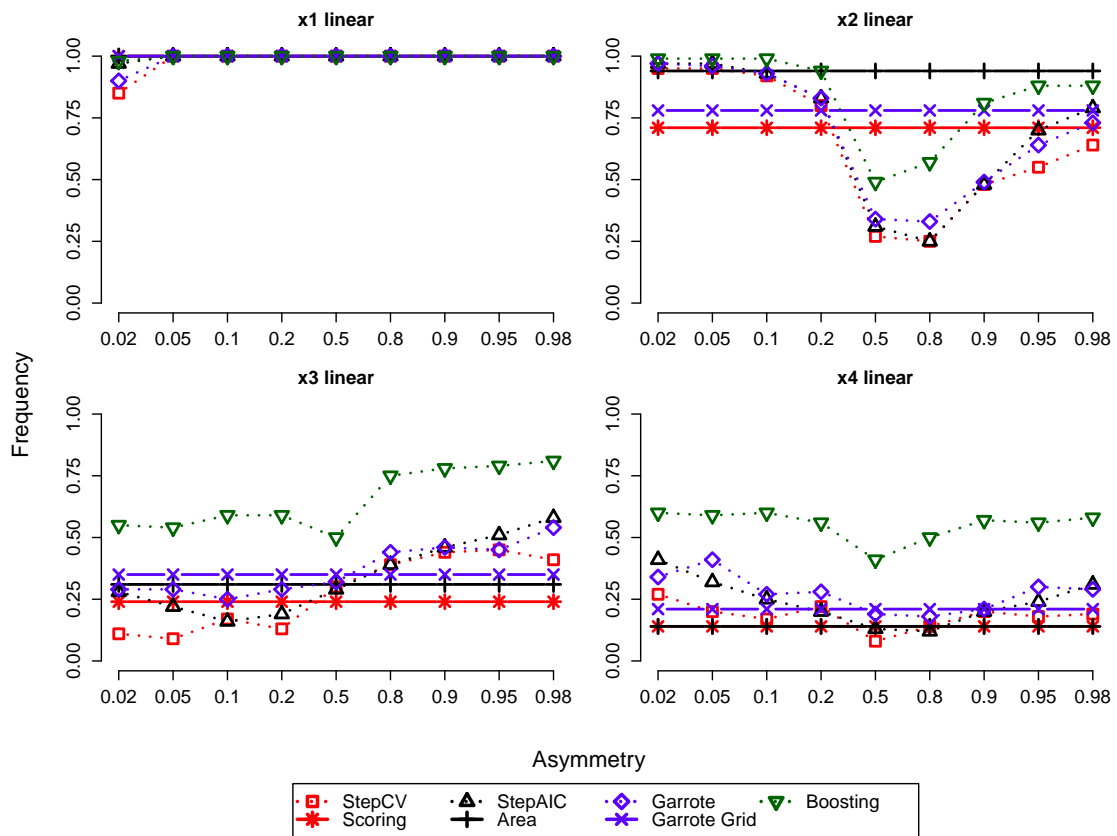


Figure A.24: Frequency of selected models for exponential design with $n=500$ and selection as linear or no effect

Part II

Application

3 Selection results for all approaches

i.) Cross-validation and scoring

Covariate	Type	0.02	0.05	0.1	0.2	0.5	0.8	0.9	0.95	0.98	grid
birth order		■	■	■	■	■	■	■	■	■	■
caesarian		■	■								
dead children		■									
household head											
household members			■	■	■	■	■	■	■	■	■
mother's education		■	■	■	■	■	■	■	■		■
partner's education						■	■	■	■	■	■
sex		■	■	■	■	■			■	■	■
bicycle											
electricity								■			
motorcycle			■	■	■	■	■	■	■	■	■
radio											
refrigerator		■	■	■	■	■	■	■	■	■	■
telephone			■	■	■	■	■	■	■	■	■
television		■	■	■	■	■	■	■	■	■	■
breastfeeding	linear										
breastfeeding	nonlinear					■	■	■	■	■	■
child's age	linear				■		■	■	■	■	■
child's age	nonlinear	■	■	■	■	■	■	■	■	■	■
mother's age	linear	■	■	■		■	■	■	■	■	■
mother's age	nonlinear		■	■	■	■			■	■	
mother's bmi	linear	■	■	■	■	■	■	■	■	■	■
mother's bmi	nonlinear										
mother's height	linear	■	■	■	■	■	■	■	■	■	■
mother's height	nonlinear	■	■								
region	GMRF	■	■	■	■	■	■	■	■	■	■

Table A.1: Selected covariates for stepwise forward selection with 10-fold cross-validation and scoring

ii.) Stepwise AIC and area under the AIC curve

Covariate	Type	0.02	0.05	0.1	0.2	0.5	0.8	0.9	0.95	0.98	grid
birth order		■	■	■	■	■	■	■	■	■	■
caesarian		■	■						■	■	
dead children											
household head											
household members		■	■	■	■	■	■	■	■	■	■
mother's education		■	■	■	■	■	■	■	■	■	■
partner's education					■	■	■	■	■	■	■
sex		■	■	■	■	■			■	■	■
bicycle											
electricity								■	■	■	
motorcycle		■	■	■	■	■	■	■	■	■	■
radio		■									
refrigerator		■	■	■	■	■	■	■	■	■	■
telephone		■	■	■	■	■	■	■	■	■	■
television		■	■	■	■	■	■	■	■	■	■
breastfeeding	linear										
breastfeeding	nonlinear	■	■	■	■	■	■	■	■	■	■
child's age	linear										
child's age	nonlinear	■	■	■	■	■	■	■	■	■	■
mother's age	linear										
mother's age	nonlinear	■	■	■	■	■	■	■	■	■	■
mother's bmi	linear	■	■	■	■	■	■				■
mother's bmi	nonlinear							■	■	■	■
mother's height	linear	■	■	■	■	■	■	■	■		■
mother's height	nonlinear	■	■	■	■		■	■	■	■	■
region	GMRF	■	■	■	■	■	■	■	■	■	■

Table A.2: Selected covariates with stepwise forward selection with AIC and area under the AIC curve

iii.) Non-negative garrote and non-negative garrote on the grid

Covariate	Type	0.02	0.05	0.1	0.2	0.5	0.8	0.9	0.95	0.98	grid
birth order		■	■	■	■	■	■	■	■	■	■
caesarian		■	■					■	■	■	
dead children											
household head											
household members		■	■	■	■	■	■	■	■	■	■
mother's education		■	■	■	■	■	■	■	■	■	■
partner's education			■	■	■	■	■	■	■	■	■
sex		■	■	■	■	■			■	■	■
bicycle											
electricity								■	■		
motorcycle			■	■	■	■	■	■	■	■	■
radio			■								
refrigerator		■	■	■	■	■	■	■	■	■	■
telephone			■	■	■	■	■	■	■	■	■
television		■	■	■	■	■	■	■	■	■	■
breastfeeding	linear					■	■	■	■	■	■
breastfeeding	nonlinear					■	■	■	■	■	■
child's age	linear		■	■	■	■	■	■	■	■	■
child's age	nonlinear	■	■	■	■	■	■	■	■	■	■
mother's age	linear	■	■	■	■	■	■	■	■	■	■
mother's age	nonlinear	■	■	■	■	■	■	■	■	■	■
mother's bmi	linear	■	■	■	■	■	■	■			■
mother's bmi	nonlinear								■	■	■
mother's height	linear	■	■	■	■	■	■	■	■	■	■
mother's height	nonlinear	■	■	■					■	■	■
region	GMRF	■	■	■	■	■	■	■	■	■	■

Table A.3: Selected covariates with non-negative garrote and non-negative garrote on the grid

iv.) Expectile Boosting

Covariate	Type	0.02	0.05	0.1	0.2	0.5	0.8	0.9	0.95	0.98
birth order		■	■	■	■	■	■	■	■	■
caesarian		■	■	■		■	■	■	■	■
dead children						■	■			
household head				■	■	■				
household members		■	■	■	■	■	■	■	■	■
mother's education		■	■	■	■	■	■	■	■	■
partner's education			■	■	■	■	■	■	■	■
sex		■	■	■	■	■			■	■
bicycle										
electricity										
motorcycle		■	■	■	■	■	■	■	■	■
radio							■			
refrigerator		■	■	■	■	■	■	■	■	■
telephone		■	■	■	■	■	■	■	■	■
television		■	■	■	■	■	■	■	■	■
breastfeeding	linear									
breastfeeding	nonlinear	■	■	■	■	■	■	■	■	■
child's age	linear			■	■	■	■	■	■	■
child's age	nonlinear	■	■	■	■	■	■	■	■	■
mother's age	linear		■	■	■	■	■	■	■	
mother's age	nonlinear	■	■	■	■	■	■	■	■	■
mother's bmi	linear					■	■	■		
mother's bmi	nonlinear	■	■	■	■	■	■	■	■	■
mother's height	linear	■	■	■	■	■	■	■	■	■
mother's height	nonlinear	■	■	■	■	■	■	■	■	■
region	GMRF	■	■	■	■	■	■	■	■	■

Table A.4: Selected covariates with Boosting and maximal value m_{stop} of 4000.

v.) Quantile Boosting

Covariate	Type	0.07	0.13	0.19	0.29	0.5	0.71	0.81	0.87	0.93
birth order		■	■	■	■	■	■	■	■	■
caesarian				■	■	■	■	■	■	■
dead children		■	■		■		■	■		
household head			■	■						
household members		■	■	■	■	■	■	■	■	■
mother's education		■	■	■	■	■	■	■	■	■
partner's education			■	■	■	■	■	■	■	■
sex		■	■	■	■	■	■			
bicycle						■				■
electricity										
motorcycle					■	■	■	■	■	■
radio							■			
refrigerator		■	■	■	■	■	■	■	■	■
telephone		■	■	■	■	■	■	■	■	■
television		■	■	■	■	■	■	■	■	■
breastfeeding	linear									
breastfeeding	nonlinear	■	■	■	■	■	■	■	■	■
child's age	linear		■	■	■	■	■	■	■	■
child's age	nonlinear	■	■	■	■	■	■	■	■	■
mother's age	linear	■	■	■	■	■	■	■		
mother's age	nonlinear	■	■	■	■	■	■	■	■	■
mother's bmi	linear			■	■	■	■	■		
mother's bmi	nonlinear	■	■	■	■	■	■	■	■	■
mother's height	linear	■	■	■	■	■	■	■	■	■
mother's height	nonlinear	■	■	■	■	■	■	■	■	■
region	GMRF	■	■	■	■	■	■	■	■	■

Table A.5: Selected covariates via Quantile Boosting and maximal value m_{stop} of 4000. The asymmetry levels are transformed under the assumption of an underlying Gaussian distribution. Then these results should coincide with the results of Expectile Boosting based on the original asymmetry levels.

vi.) Selection based on confidence intervals

Covariate	Type	0.02	0.05	0.1	0.2	0.5	0.8	0.9	0.95	0.98
birth order		■	■	■	■	■	■	■	■	■
caesarian		■								
dead children										
household head										
household members		■	■	■	■	■	■	■	■	■
mother's education		■	■	■	■	■	■	■	■	■
partner's education						■	■	■	■	■
sex		■	■	■	■					
bicycle										
electricity										
motorcycle						■	■	■	■	■
radio										
refrigerator		■	■	■	■	■	■	■	■	■
telephone			■	■	■	■	■	■	■	■
television		■	■	■	■	■	■	■	■	■
breastfeeding	linear									
breastfeeding	nonlinear pointwise				■	■	■	■	■	■
breastfeeding	nonlinear SCB									
child's age	linear				■	■	■	■	■	■
child's age	nonlinear pointwise	■	■	■	■	■	■	■	■	■
child's age	nonlinear SCB	■	■	■	■	■	■	■	■	■
mother's age	linear				■	■	■	■	■	■
mother's age	nonlinear pointwise	■	■	■	■	■	■	■	■	■
mother's age	nonlinear SCB									
mother's bmi	linear	■	■	■	■	■				
mother's bmi	nonlinear pointwise							■	■	■
mother's bmi	nonlinear SCB									
mother's height	linear	■	■	■	■	■	■	■	■	■
mother's height	nonlinear pointwise		■	■						
mother's height	nonlinear SCB									
region	GMRF	■	■	■	■	■	■	■	■	■

Table A.6: Selected covariates based on bootstrap confidence intervals, with 2000 bootstrap replications. The nonlinear functions are treated with simultaneous (SCB) and pointwise confidence intervals.

vii.) Standard scoring vs. weighted scoring

Covariate	Type	standard	weighted
birth order		■	■
caesarian			
dead children			
household head			
household members		■	■
mother's education		■	■
partner's education		■	■
sex		■	■
bicycle			
electricity			
motorcycle		■	■
radio			
refrigerator		■	■
telephone		■	■
television		■	■
breastfeeding	linear		
breastfeeding	nonlinear	■	■
child's age	linear	■	■
child's age	nonlinear	■	■
mother's age	linear	■	■
mother's age	nonlinear		■
mother's bmi	linear	■	■
mother's bmi	nonlinear		
mother's height	linear	■	■
mother's height	nonlinear		
region	GMRF	■	■

Table A.7: Comparison of selected covariates via 10-fold scoring, unweighted and with a weight of 10 for all asymmetries smaller than 0.11

viii.) Stepwise backward CV after scoring

Covariate	Type	0.02	0.05	0.1	0.2	0.5	0.8	0.9	0.95	0.98
birth order		■	■	■	■	■	■	■	■	■
household members			■	■	■	■	■	■	■	■
mother's education		■	■	■	■	■	■	■	■	■
partner's education						■	■	■	■	■
sex		■	■	■	■			■	■	■
motorcycle			■	■	■	■	■	■	■	■
refrigerator		■	■	■	■	■	■	■	■	■
telephone			■	■	■	■	■	■	■	■
television		■	■	■	■	■	■	■	■	■
breastfeeding	nonlinear		■			■	■	■	■	■
child's age	linear		■	■	■		■	■	■	■
child's age	nonlinear	■	■	■	■	■	■	■	■	■
mother's age	linear	■	■	■	■	■	■	■	■	■
mother's bmi	linear	■	■	■	■	■	■	■	■	■
mother's height	linear	■	■	■	■	■	■	■	■	■
region	GMRF	■	■	■	■	■	■	■	■	■

Table A.8: Selected covariates with backward 10-fold CV after scoring

4 Estimated regional effects for all used asymmetries

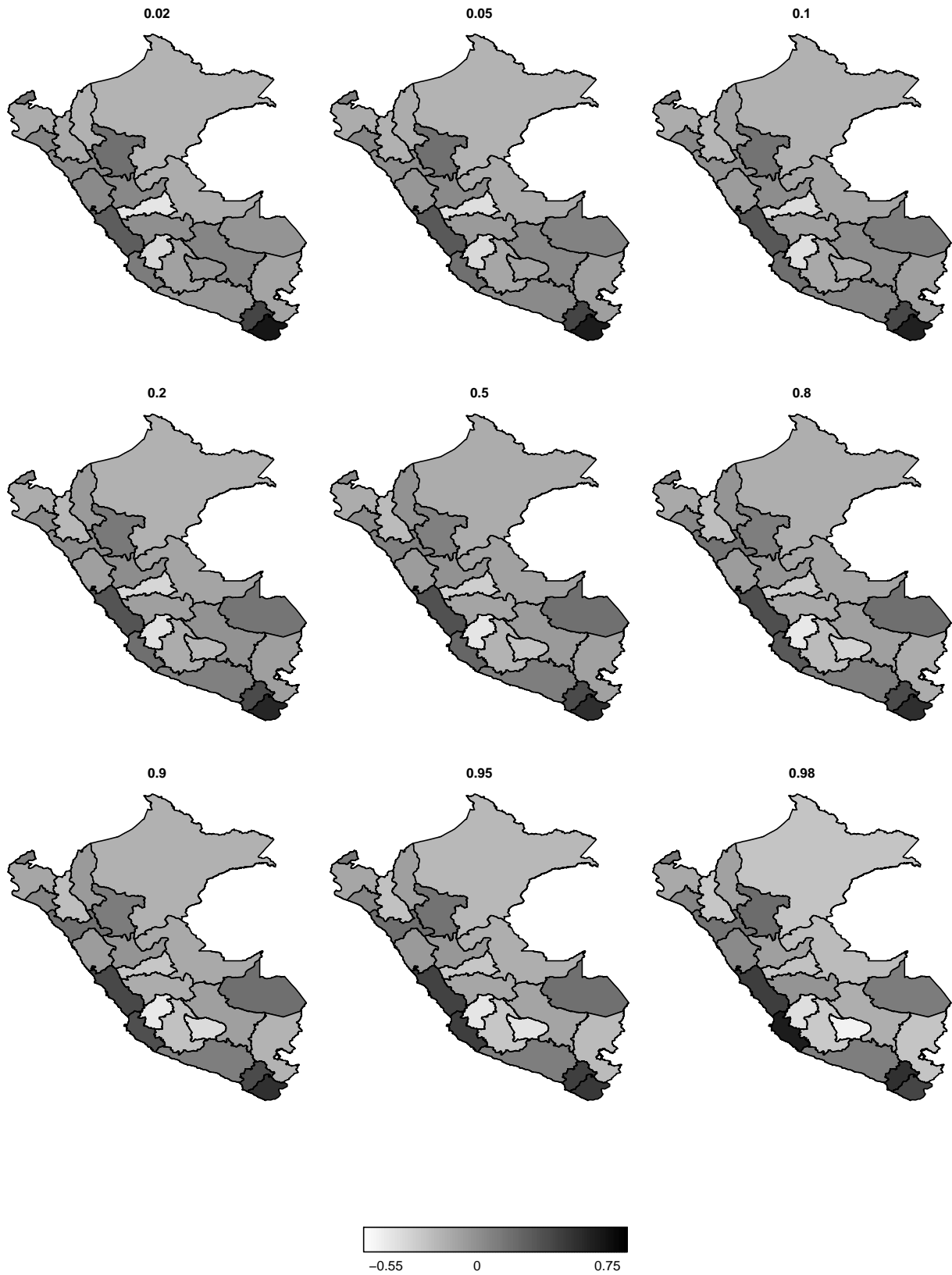


Figure A.25: Estimated coefficients of regional information for the best model selected via scoring