

Firth's penalized likelihood logistic regression: accurate effect estimates AND predictions?

Angelika Geroldinger, Rainer Puhr, Mariana Nold, Lara Lusa, Georg Heinze

15.7.2016

XXVIIIth International Biometric Conference



Example: Bias in logistic regression

Consider a model containing only intercept, no regressors:

$$\text{logit}(P(Y = 1)) = \beta_0.$$

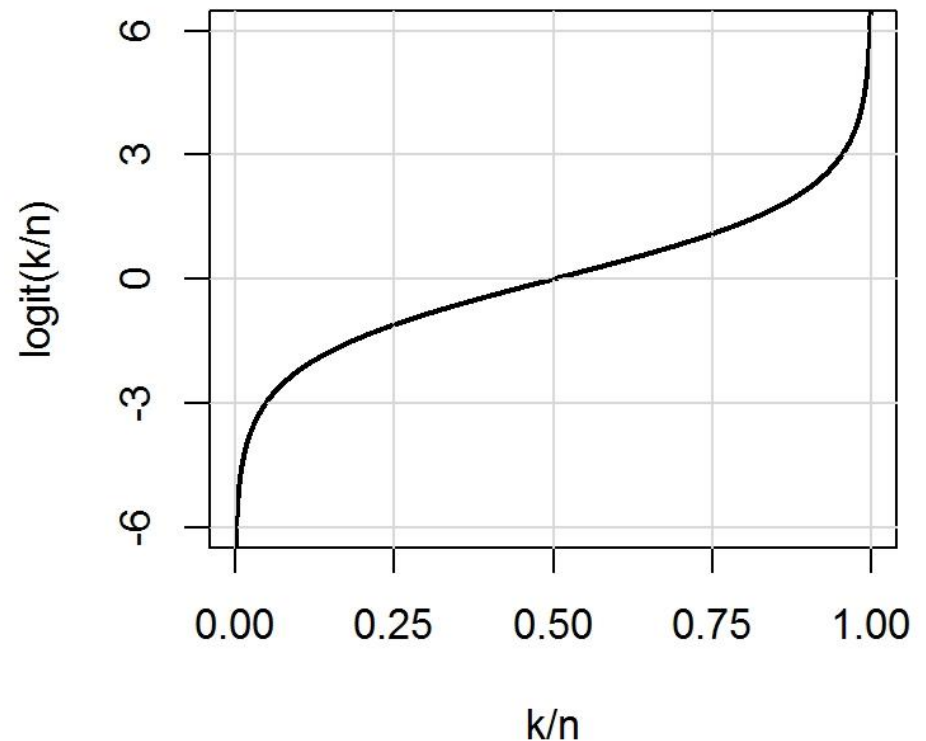
With n observations, k events, the ML estimator of β_0 is given by:

$$\widehat{\beta}_0 = \text{logit}(k/n).$$

Since k/n is unbiased,

$\widehat{\beta}_0$ is biased!

(If $\widehat{\beta}_0$ was unbiased,
 $\text{expit}(\widehat{\beta}_0)$ would be biased!)



Firth type penalization

In exponential family models with canonical parametrization the **Firth-type penalized likelihood** is given by

$$L^*(\beta) = L(\beta) \det(I(\beta))^{1/2},$$

where $I(\beta)$ is the Fisher information matrix and $L(\beta)$ is the likelihood.

Firth-type penalization

- **removes the first-order bias** of the ML-estimates of β ,
- is **bias-preventive** rather than corrective,
- is available in **Software** packages such as SAS, R, Stata...

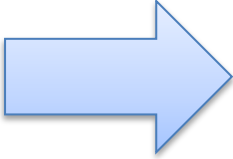
Firth type logistic regression

In logistic regression, the penalized likelihood is given by

$$L^*(\beta) = L(\beta) \det(X^t W X)^{1/2}, \text{ with}$$

$$\begin{aligned} W &= \text{diag}(\text{expit}(X_i \beta)(1 - \text{expit}(X_i \beta))) \\ &= \text{diag}(\pi_i(1 - \pi_i)). \end{aligned}$$

- Firth-type estimates always exist.

 W is maximised at $\pi_i = \frac{1}{2}$, i.e. at $\beta = 0$, thus

- predictions are usually pulled towards $\frac{1}{2}$,
- coefficients towards zero.

Firth type logistic regression (FL)

For logistic regression with one binary regressor,
Firth's bias correction amounts to adding 1/2 to each cell:

original			augmented			
	A	B		A	B	
0	44	4	Firth-type penalization →	0	44.5	4.5
1	1	1		1	1.5	1.5

$$\text{event rate} = \frac{2}{50} = 0.04$$

$$\text{OR}_{B \text{ vs } A} = 11$$

$$\text{event rate} = \frac{3}{52} \sim 0.058$$

$$\text{OR}_{B \text{ vs } A} = 9.89$$

$$\text{av. pred. prob.} = 0.054$$

FLAC

Split the augmented data into the original and pseudo data:

augmented			original			pseudo				
	A	B		A	B		A	B		
0	44.5	4.5	→	0	44	4	+	0	0.5	0.5
1	1.5	1.5		1	1	1		1	0.5	0.5

Define **F**irth type **L**ogistic regression with **A**dditional **C**ovariate as the stratified analysis of the original and pseudo data:

$$OR_{BvsA} = 6.63$$

av. pred. prob. = 0.04 = observed proportion of events!

FLAC

In the general case (idea):

One can show, that Firth-type penalization is equivalent to ML estimation of augmented data.

FLAC estimates can be obtained by the following steps:

- 1) Define an indicator variable discriminating between original and pseudo data.
- 2) Apply ML on the augmented data including the indicator.



unbiased pred. probabilities

FLIC

Firth type Logistic regression with Intercept Correction:

Modify the intercept in Firth-type estimates such that the average pred. prob. becomes equal to the observed proportion of events.

- ✓ unbiased pred. probabilities
- effect estimates are the same as in Firth type logistic regression

Other methods for accurate pred.

In our simulation study, we compared FLIC and FLAC to the following methods:

- weakened Firth-type penalization, with $L(\beta)^* = L(\beta) \det(X^t W X)^\tau$, $\tau < 1/2$, (WF)
- ridge regression, (RR)
- penalization by log-F(1,1) priors, (LF)
- penalization by Cauchy priors with scale parameter=2.5. (CP)

log-F(1,1) priors (LF)

Penalizing by log-F(1,1) prior gives $L(\beta)^* = L(\beta) \cdot \prod \frac{e^{\frac{\beta_j}{2}}}{1+e^{\beta_j}}$.

This amounts to the following modification of the data set:

	x1	x2	y
1	*	*	*
1	*	*	*
1	*	*	*
1	*	*	*
1	*	*	*
1	*	*	*
1	*	*	*

→

	x1	x2	y
1	*	*	*
1	*	*	*
1	*	*	*
1	*	*	*
1	*	*	*
1	*	*	*
1	*	*	*
0	1	0	0
0	1	0	1
0	0	1	0
0	0	1	1

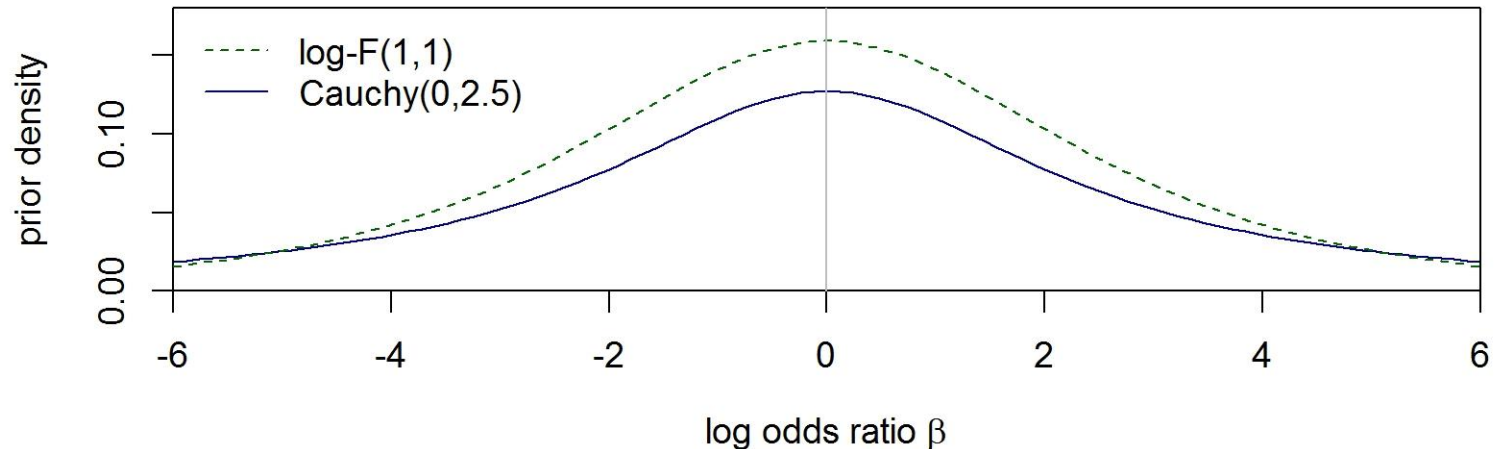
} of weight ½

We follow Greenland and Mansournia, 2015:

- no penalization of the intercept,
- no scaling of variables.

Cauchy priors (CP)

Cauchy priors (scale=2.5) have heavier tails than log-F(1,1)-priors:



We follow Gelman et al., 2008:

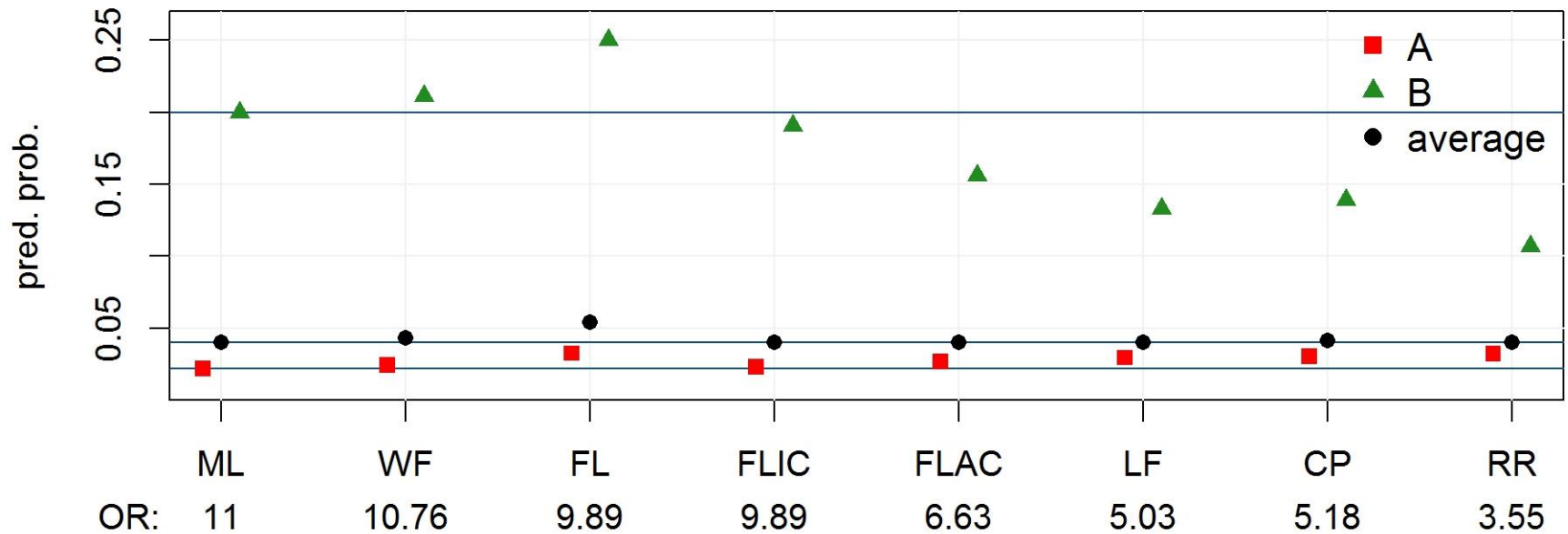
- all variables are centered,
- binary variables are coded to have a range of 1,
- all other variables are scaled to have standard deviation 0.5,
- the intercept is penalized by Cauchy(0,10).

This is implemented in the function `bayesglm` in the R-package `arm`.

Revisiting the toy example

	A	B
0	44	4
1	1	1

The different methods give:



unbiased
pred. prob.:



Simulation study: the set-up

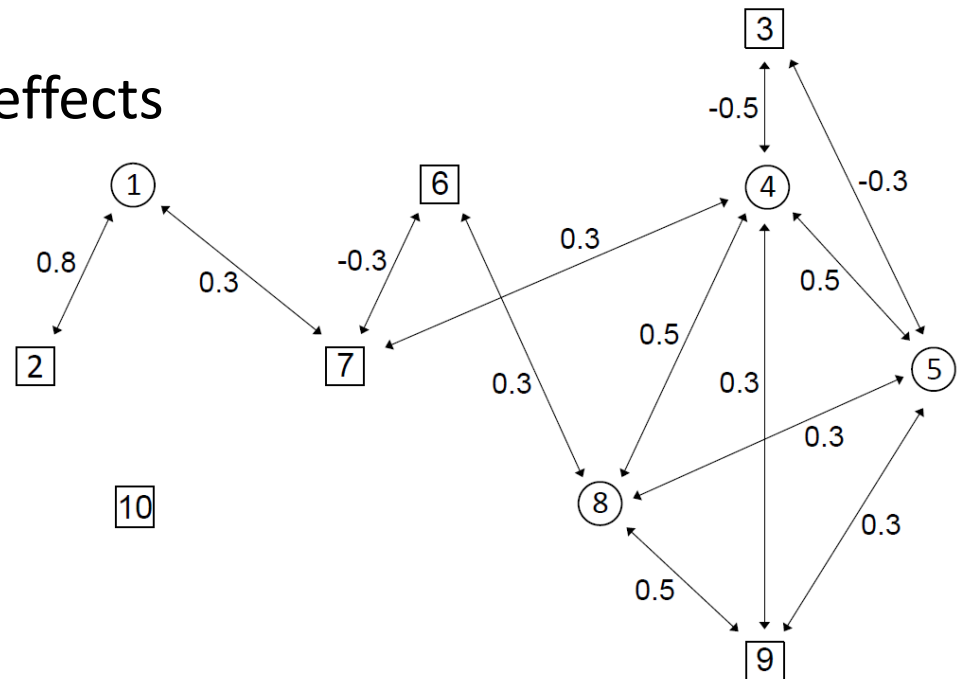
We investigated the performance of FLIC and FLAC, simulating 1000 data sets for 45 scenarios with:

- 500, 1000 or 1400 observations,
- event rates of 1%, 2%, 5% or 10%
- 10 covariables (6 cat., 4 cont.),
see Binder et al., 2011
- none, moderate and strong effects
of positive and mixed signs

Main evaluation criteria:

bias and RMSE of

- predictions and
- effect estimates



Average predicted probability

For the scenarios with small effect size:

N	method	rel.bias				rel.RMSE			
		exp. event rate				exp. event rate			
		0.01	0.02	0.05	0.1	0.01	0.02	0.05	0.1
500	WF			3.7	1.6			3.8	1.6
	FL			18.2	7.8			18.7	7.9
	CP			0.2	0.1			0.2	0.1
1400	WF		3.7	1.3	0.6		3.8	1.3	0.6
	FL		18.5	6.6	2.8		19.0	6.7	2.8
	CP		0.2	0.1	0.0		0.3	0.1	0.0
3000	WF	3.6	1.7	0.6	0.3	3.7	1.7	0.6	0.3
	FL	17.9	8.6	3.1	1.3	18.3	8.6	3.1	1.3
	CP	0.3	0.1	0.0	0.0	0.3	0.1	0.0	0.0

All other methods (ML, FLIC, FLAC, LF and RR) yield average pred. prob. equal to the proportion of events.

Average predicted probability

For the scenarios with coefficients of mixed signs and small effect size:

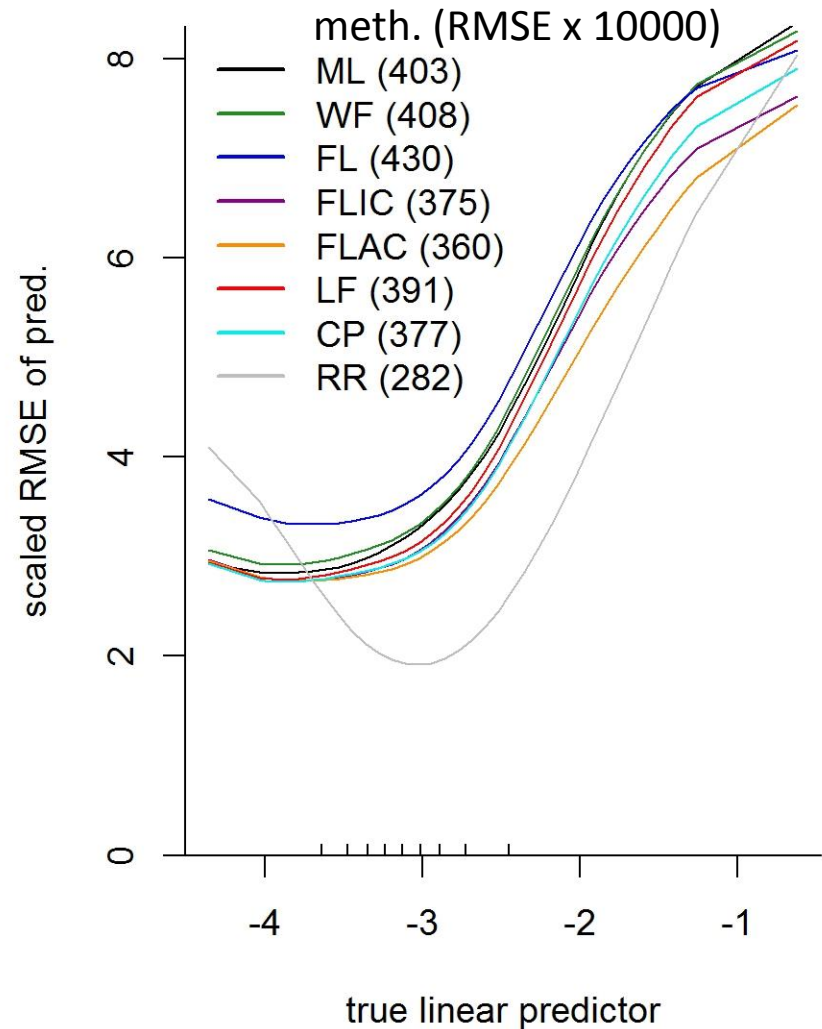
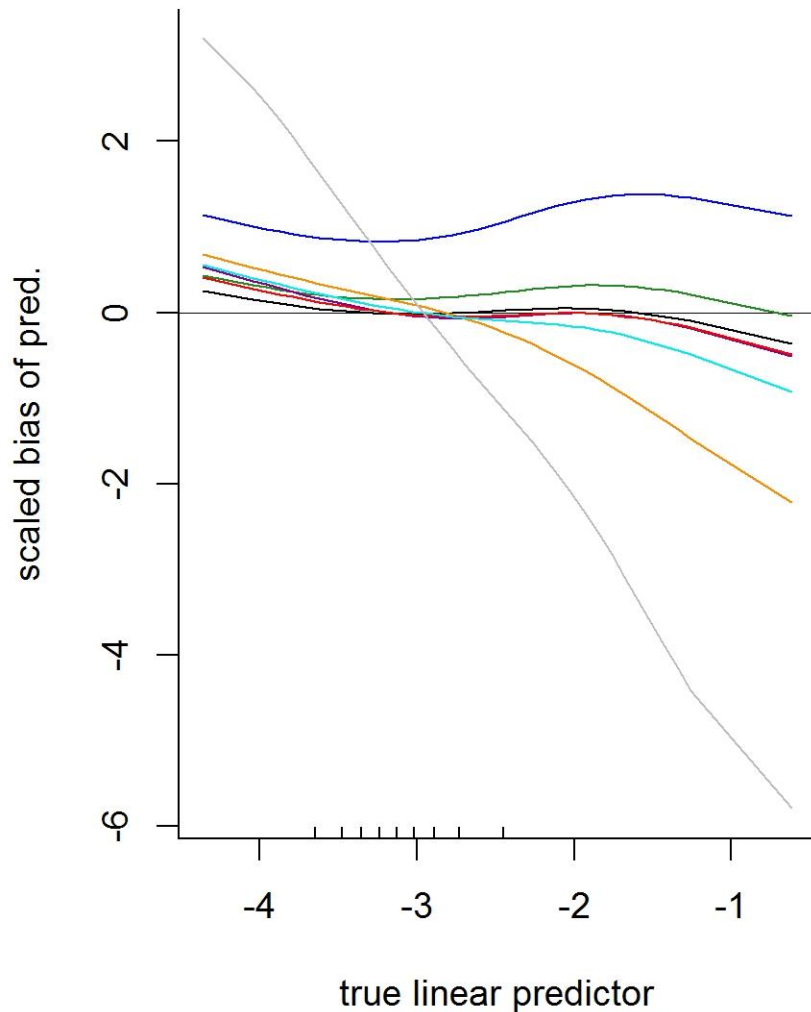
N	method	rel.bias				rel.RMSE			
		exp. event rate				exp. event rate			
		0.01	0.02	0.05	0.1	0.01	0.02	0.05	0.1
500	WF			3.7	1.6			3.8	1.6
	FL			18.2	7.8			18.7	7.9
	CP			0.2	0.1			0.2	0.1
1400	WF		3.7	1.3	0.6		1.3	0.6	0.3
	FL		18.5	6.7	2.8		6.7	2.8	1.3
	CP		0.2	0.1	0.0		0.1	0.0	0.0
3000	WF	3.6	1.7	0.6	0.3	3.7	1.7	0.6	0.3
	FL	17.9	8.6	3.1	1.3	18.3	8.6	3.1	1.3
	CP	0.3	0.1	0.0	0.0	0.3	0.1	0.0	0.0

Next, we have a closer look at this scenario...

All other methods (ML, FLIC, FLAC, LF and RR) yield average pred. prob. equal to the proportion of events.

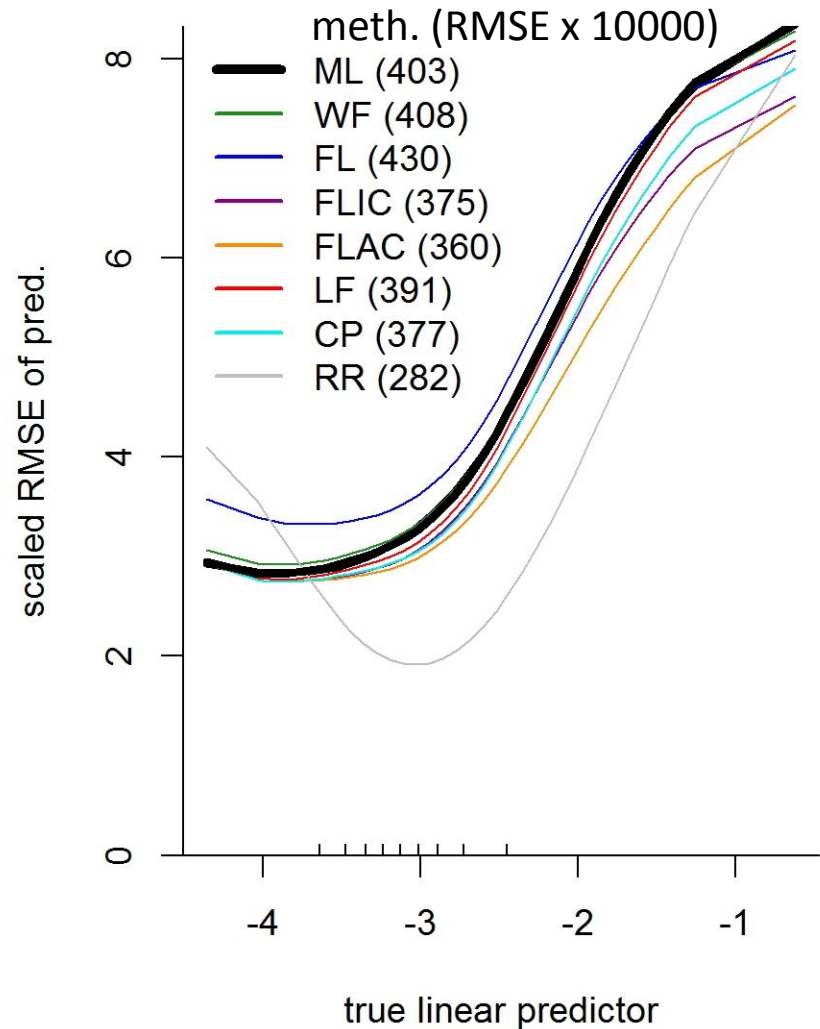
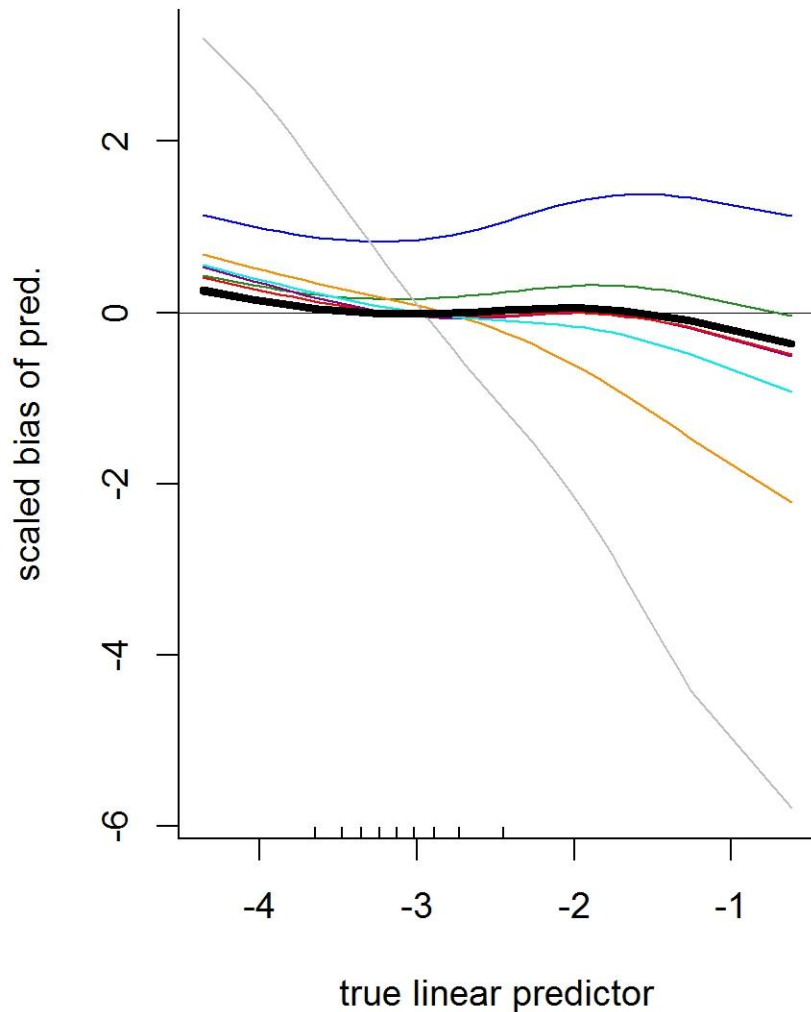
Predictions by true lin. pred.

sample size=500, prop. of events= 5%, small effect size



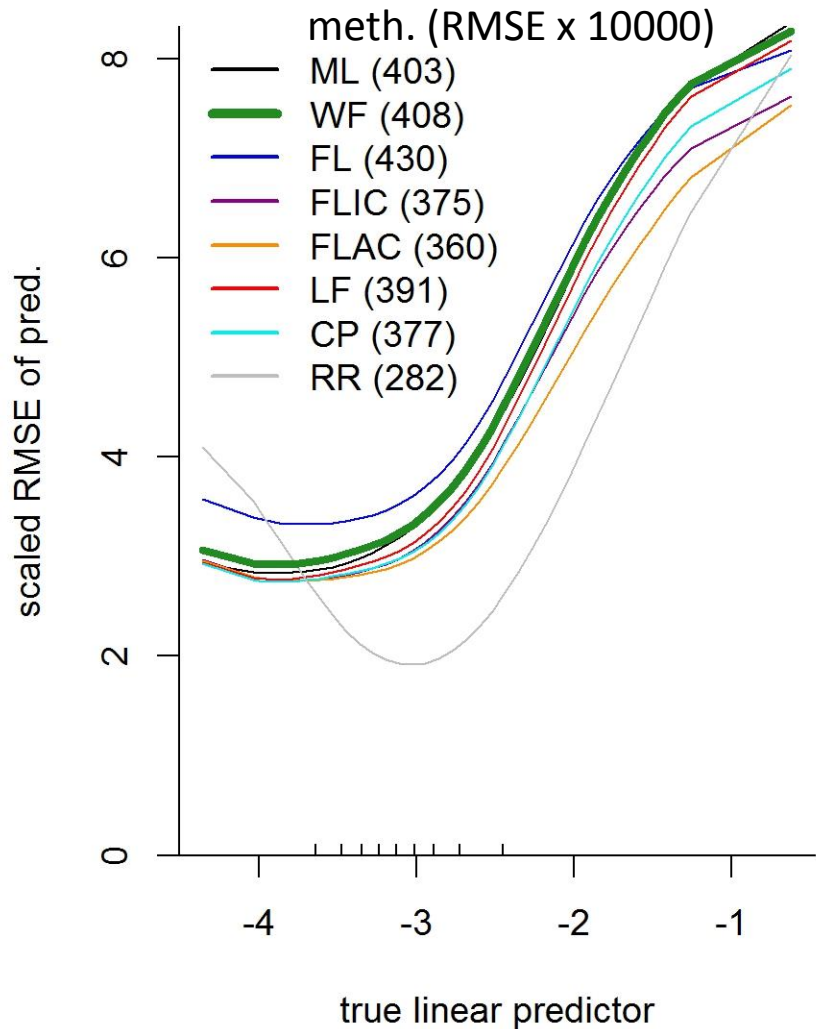
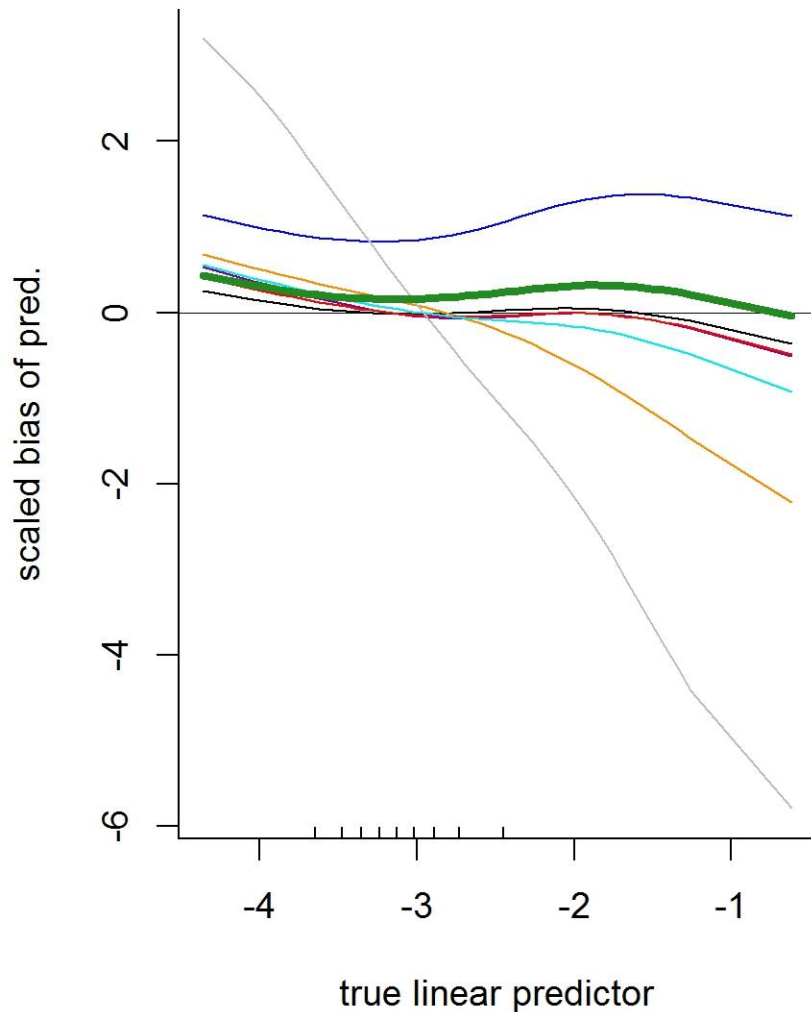
Predictions by true lin. pred.

sample size=500, prop. of events= 5%, small effect size



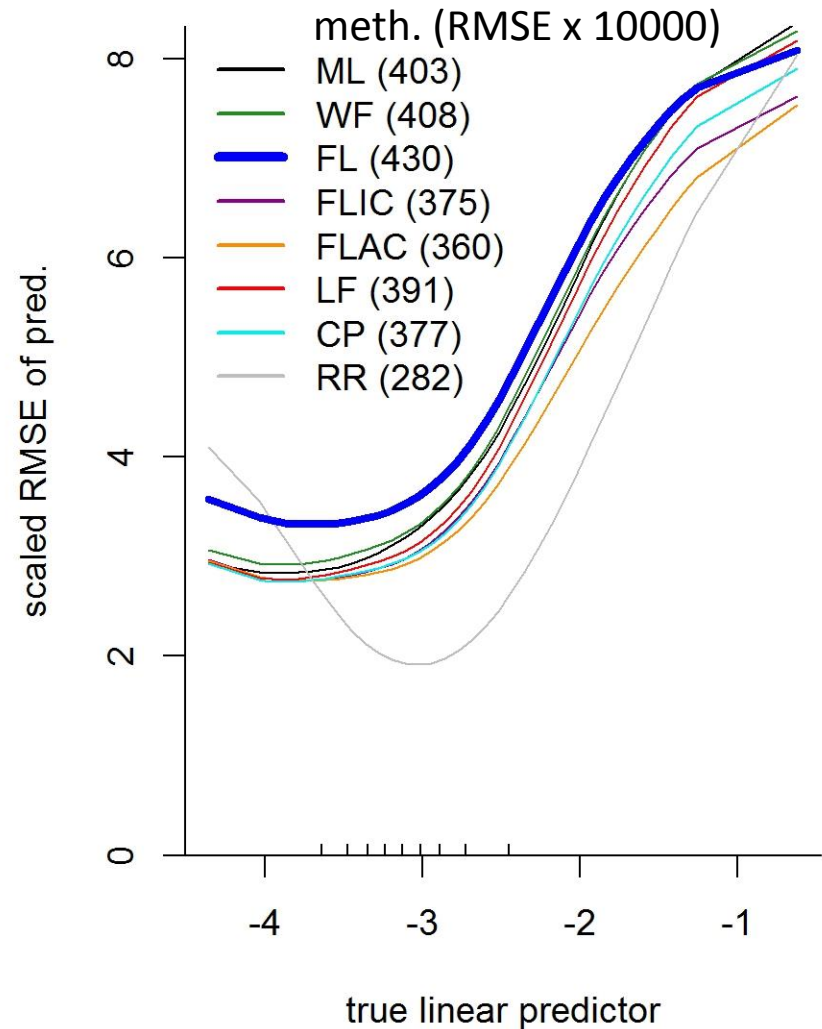
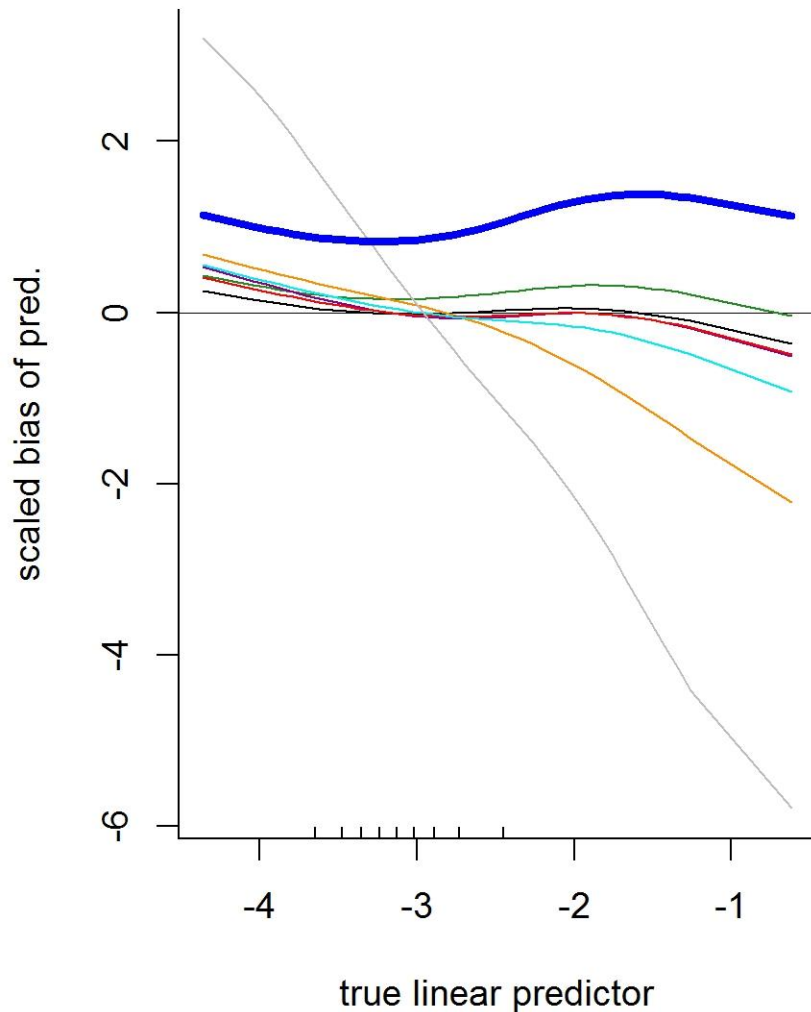
Predictions by true lin. pred.

sample size=500, prop. of events= 5%, small effect size



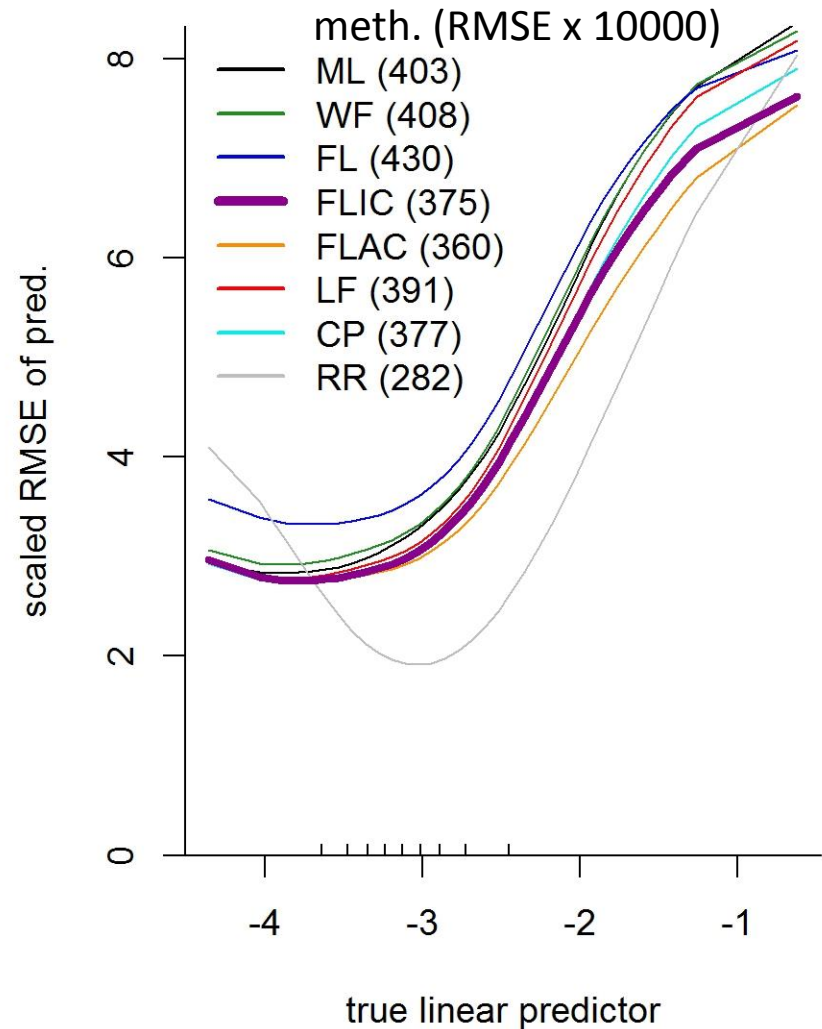
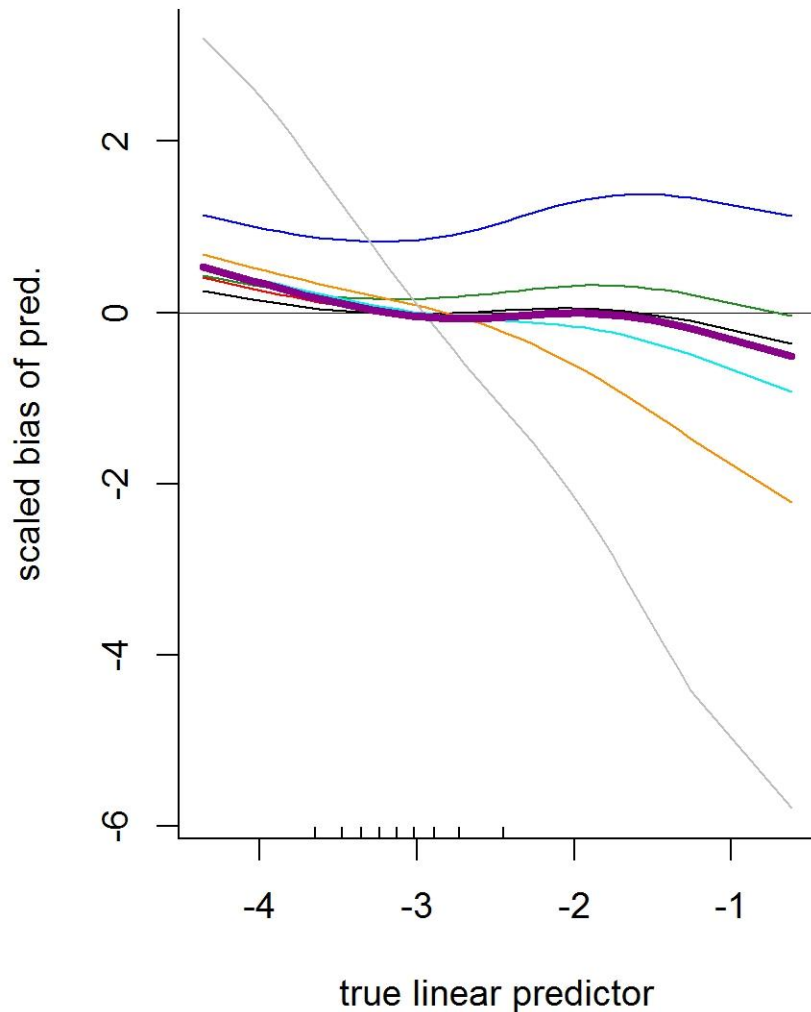
Predictions by true lin. pred.

sample size=500, prop. of events= 5%, small effect size



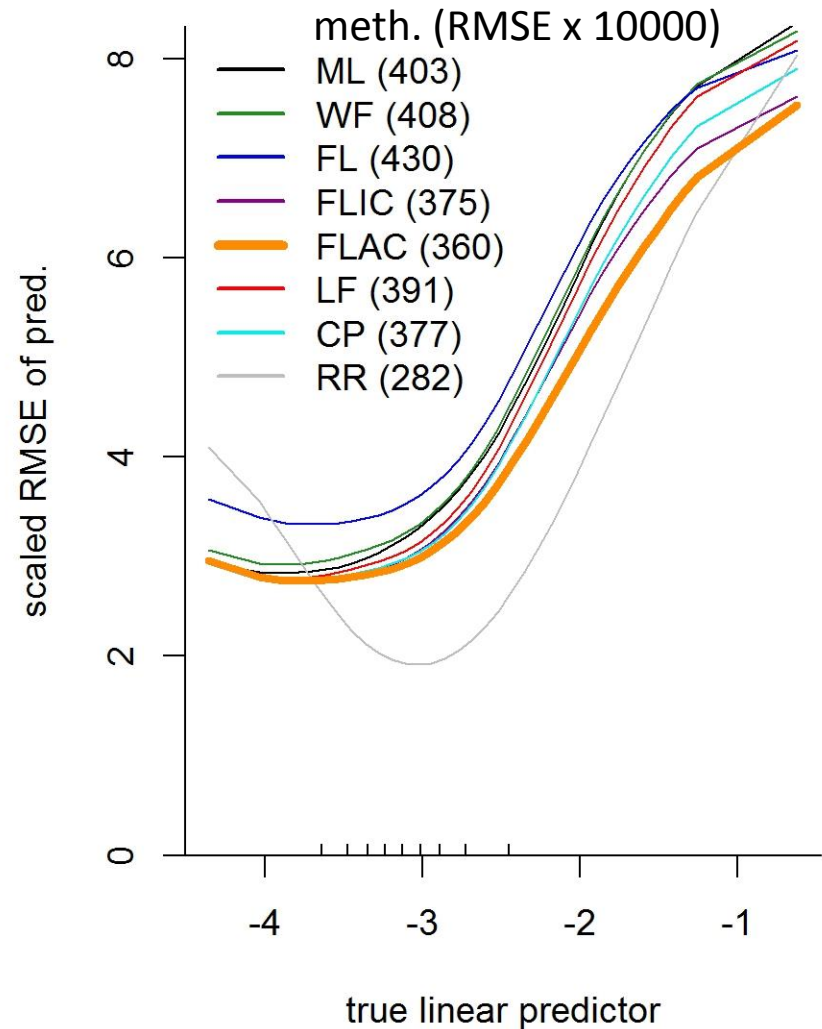
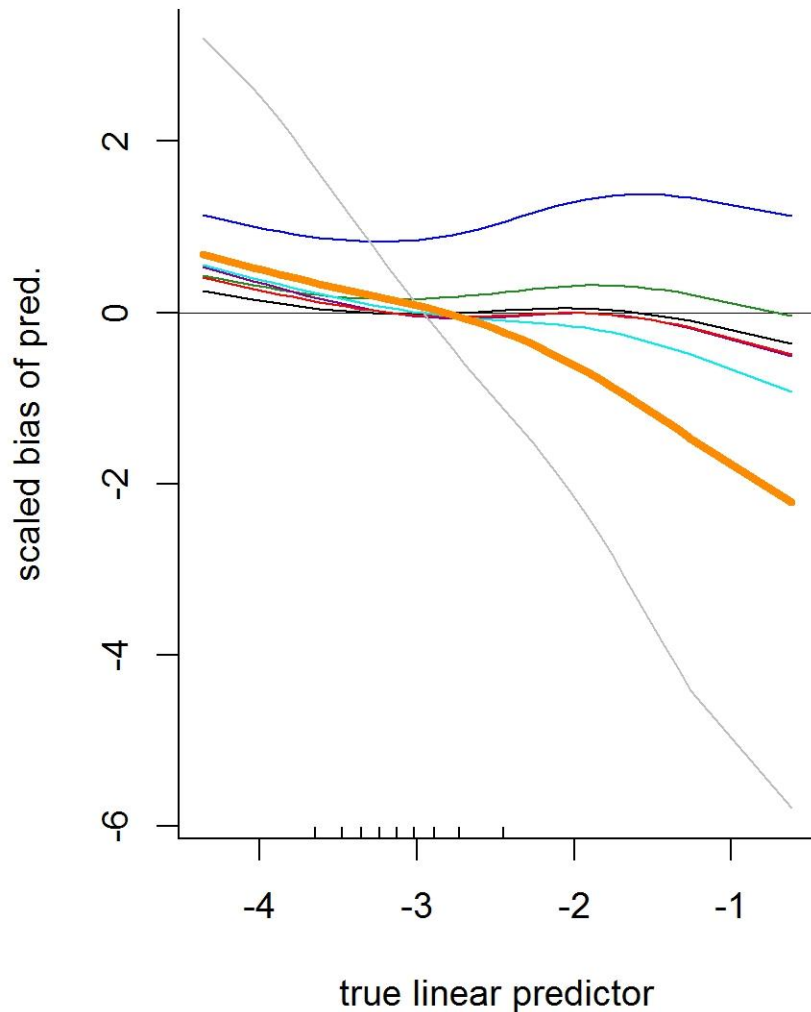
Predictions by true lin. pred.

sample size=500, prop. of events= 5%, small effect size



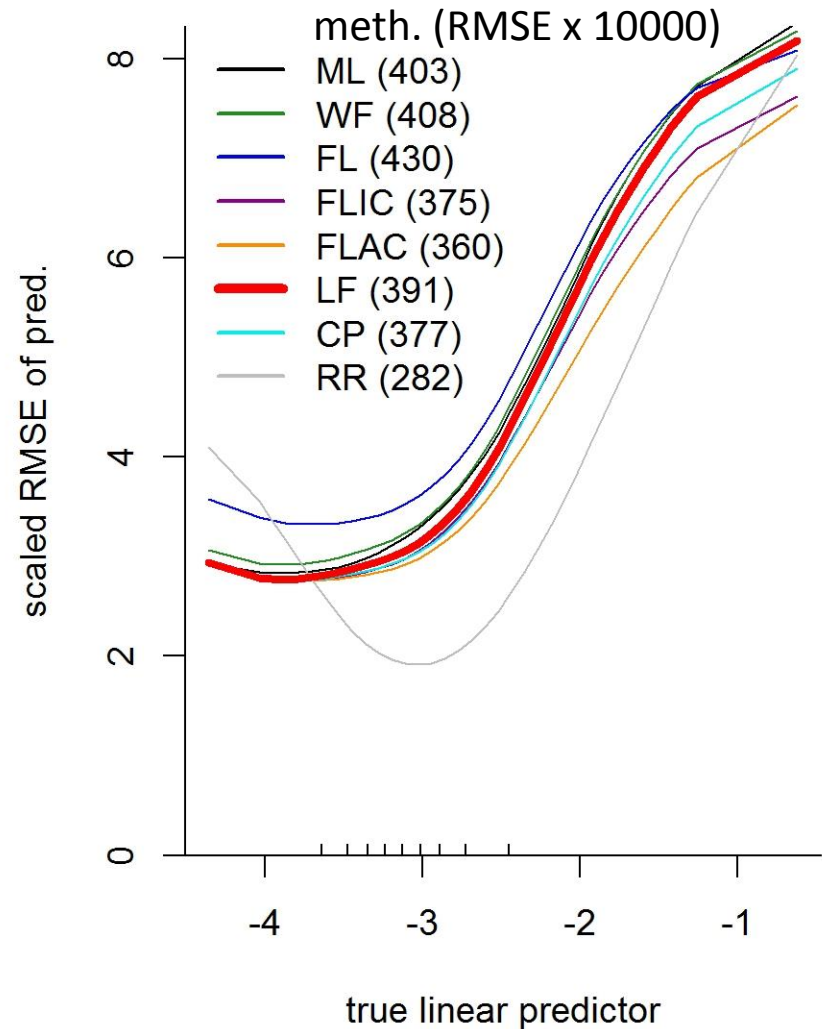
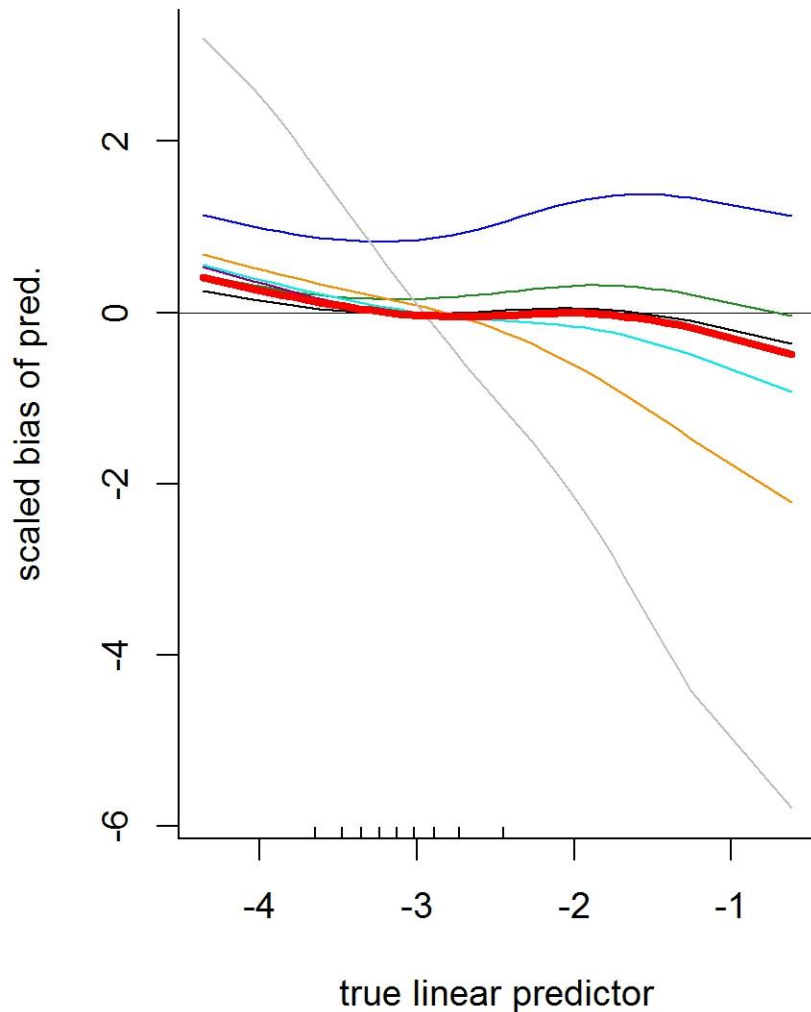
Predictions by true lin. pred.

sample size=500, prop. of events= 5%, small effect size



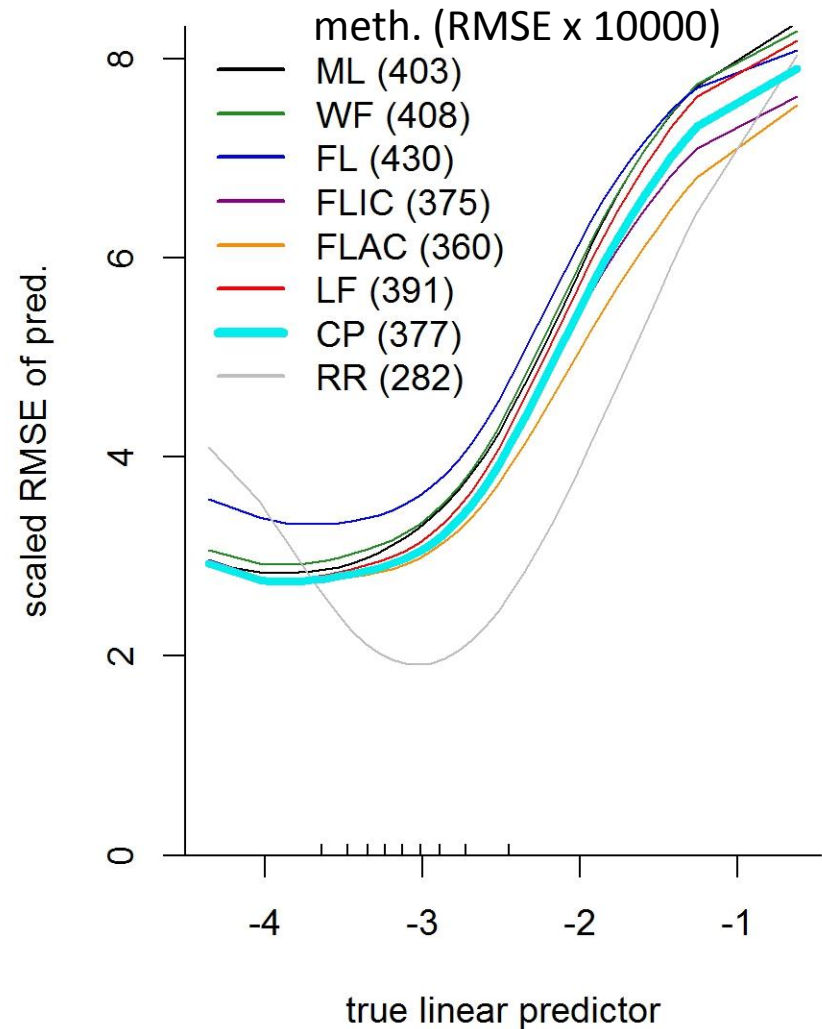
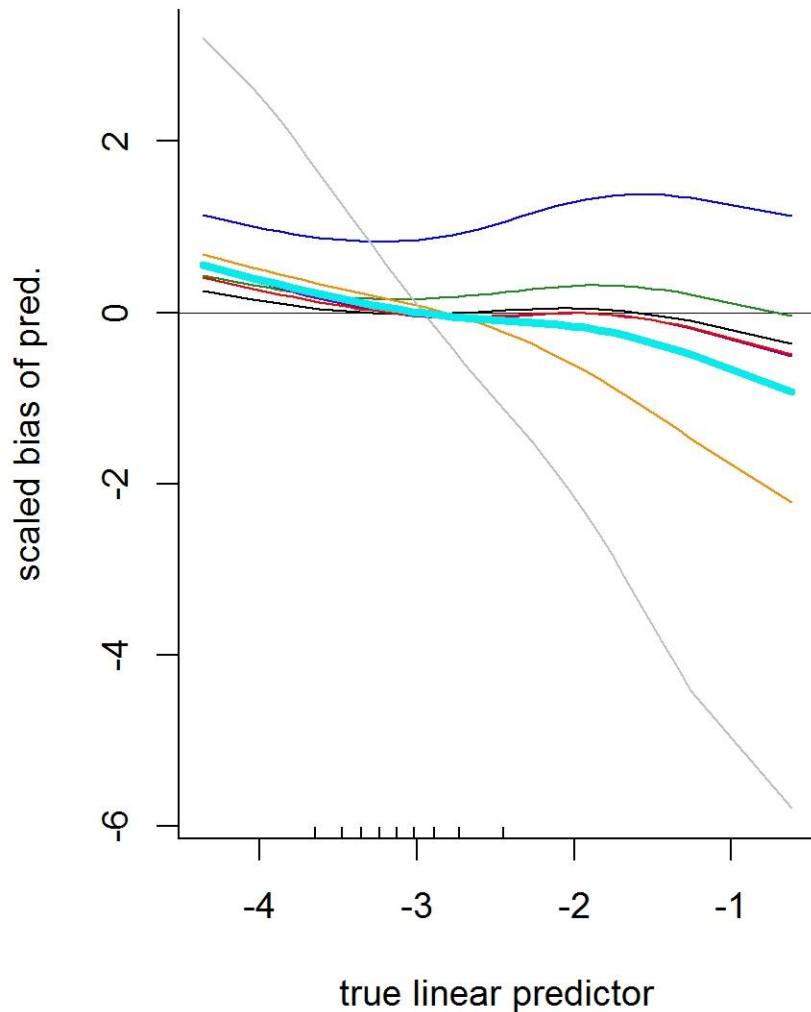
Predictions by true lin. pred.

sample size=500, prop. of events= 5%, small effect size



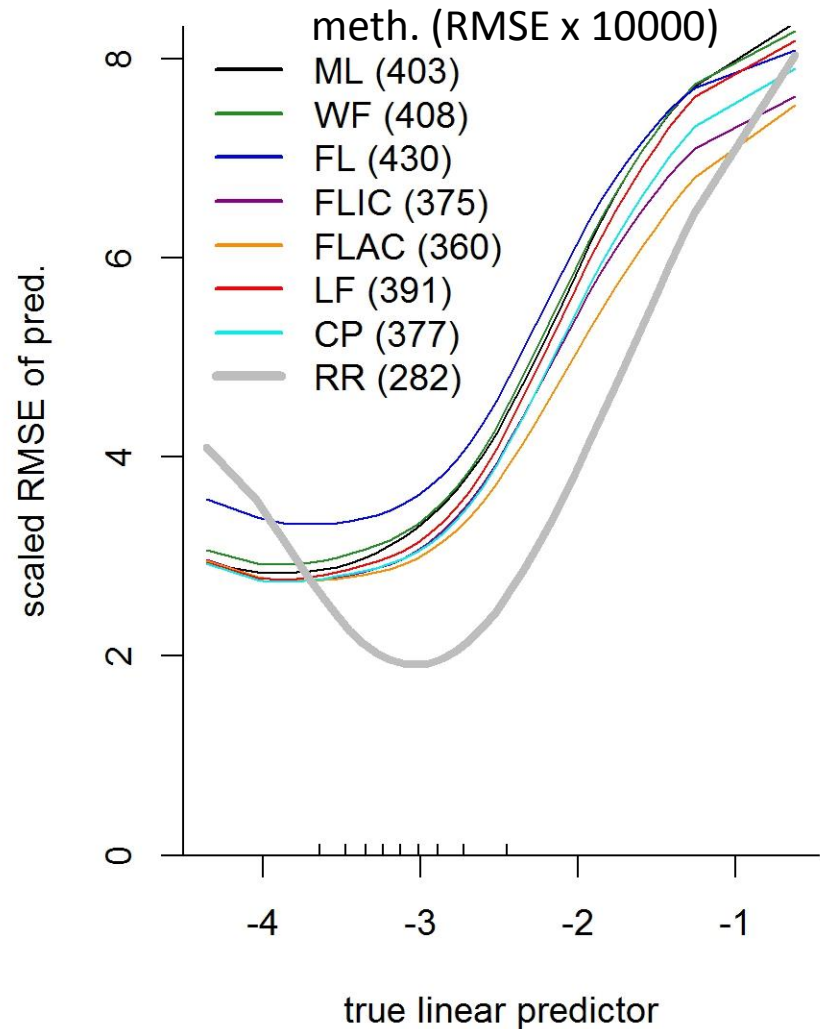
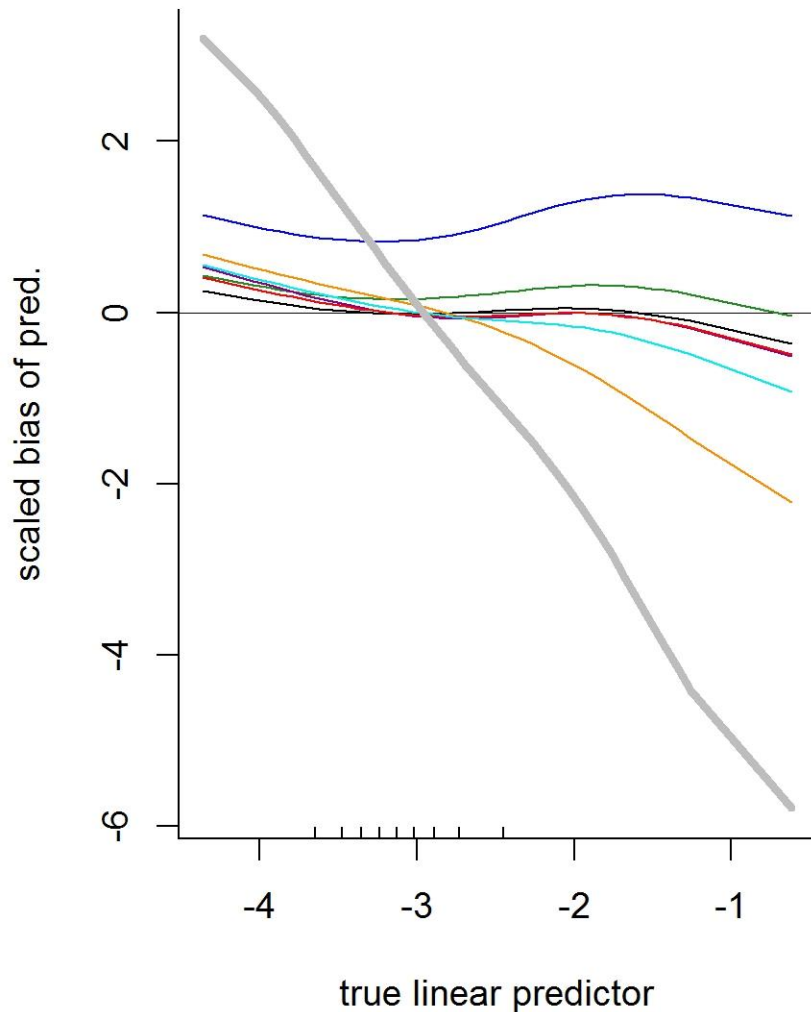
Predictions by true lin. pred.

sample size=500, prop. of events= 5%, small effect size



Predictions by true lin. pred.

sample size=500, prop. of events= 5%, small effect size



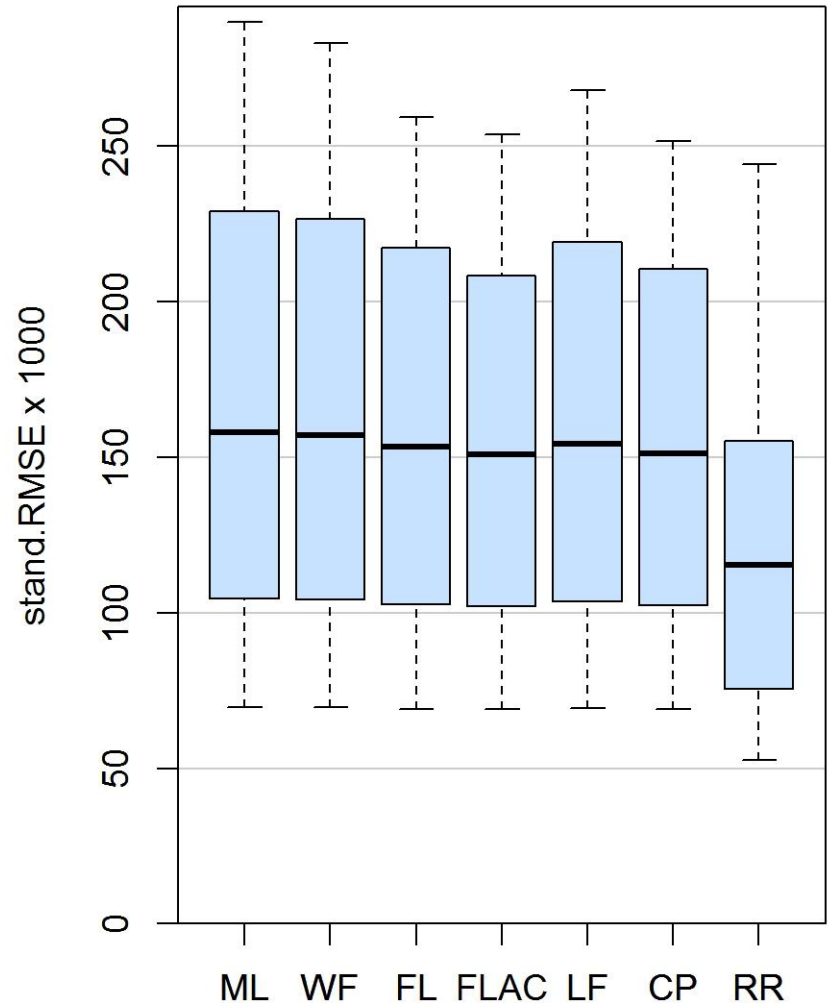
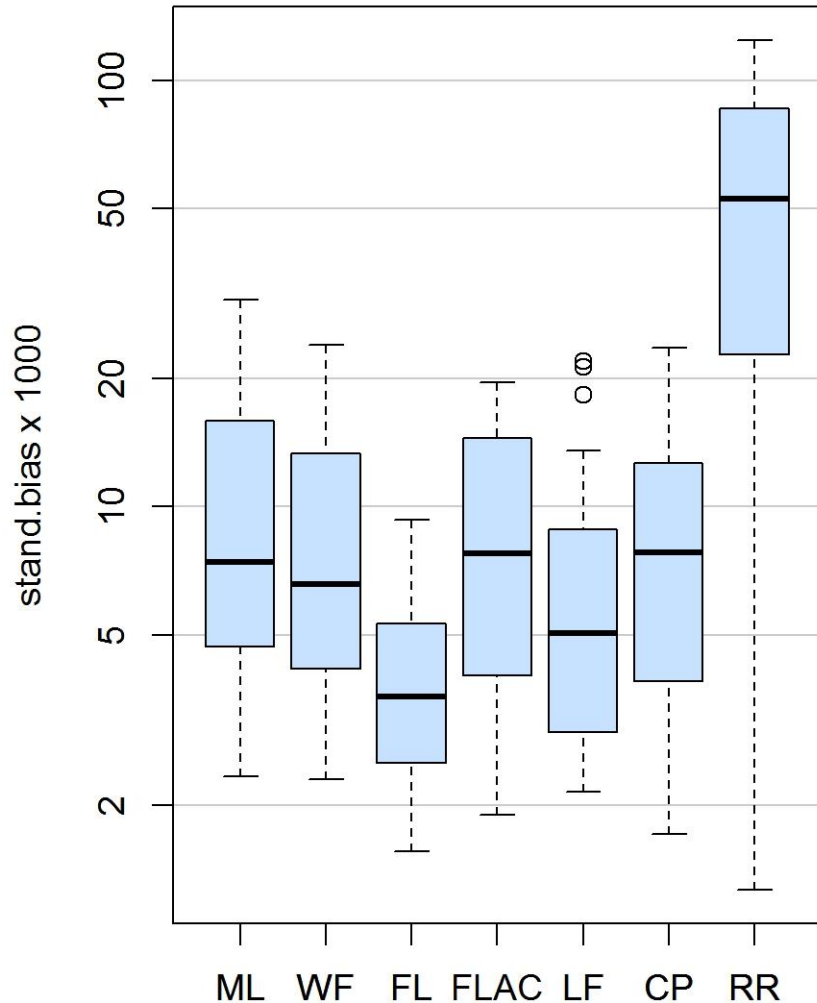
Coefficients

Absolute bias and RMSE of standardized coefficients, averaged over all 10 covariables excluding the intercept:

N	E(y)	method	bias ($\times 1000$) effect size			RMSE ($\times 1000$) effect size		
			0	0,5	1	0	0,5	1
500	0.05	ML	23	17	29	277	266	288
		WF	19	14	21	272	261	281
		FL/FLIC	7	5	9	253	244	259
		FLAC	17	16	16	239	235	252
		LF	22	10	12	265	252	266
		CP	18	14	24	245	238	251
		RR	3	109	124	78	166	244

Coefficients

Similar patterns can be observed across all 45 scenarios:



Conclusions

- For rare events, FL-predictions are severely biased (relative bias of up to 20% in our simulations).
- Both, FLIC and FLAC improved on predictions by FL, with identical effect estimates or effect estimates of lower RMSE.
- RR outperformed all other methods with respect to RMSE of coefficients and predictions, but introduces bias towards 0. Confidence intervals?
- LF performed slightly worse than CP. (Due to data preprocessing?)

Based on our simulations, if one is interested in effect estimates and predictions, we recommend to use

- RR (whenever confidence intervals are not needed)
- FLAC as a compromise between optimization of bias and RMSE.

Literature

- Binder H, Sauerbrei W and Royston P. Multivariable Model-Building with Continuous Covariates: Performance Measures and Simulation Design 2011. Technical Report FDM-Preprint 105, University of Freiburg, Germany.
- Elgmati E, Fiaccone RL, Henderson R and Matthews JNS. Penalised logistic regression and dynamic prediction for discrete-time recurrent event data. *Lifetime Data Analysis* 2015; 12(4): 542-560.
- Firth D. Bias reduction of maximum likelihood estimates. *Biometrika* 1993; 80(1): 27-38.
- Gelman A, Jakulin A, Pittau M and Su YS. A Weakly Informative Default Prior Distribution for Logistic and Other Regression Models. *The Annals of Applied Statistics* 2008; 2(4):1360–1383.
- Greenland S and Mansournia M. Penalization, Bias Reduction, and Default Priors in Logistic and Related Categorical and Survival regressions. *Statistics in Medicine* 2015; 34(23):3133-3143.
- Heinze G and Schemper M. A solution to the problem of separation in logistic regression. *Statistics in Medicine* 2002; 21(16): 2409-2419.
- Kosmidis I. Bias in parametric estimation: reduction and useful side-effects. *WIRE Computational Statistics* 2014; 6(3): 185-196.
- Puhr R, Heinze G, Nold M, Lusa L and Geroldinger A. Predicting rare events with penalized logistic regression. Work in progress.

Firth type penalization

In exponential family models with canonical parametrization the **Firth-type penalized likelihood** is given by **Jeffrey's invariant prior**

$$L^*(\beta) = L(\beta) \det(I(\beta))^{1/2},$$

where $I(\beta)$ is the Fisher information matrix and $L(\beta)$ is the likelihood.

Firth-type penalization

- **removes the first-order bias** of the ML-estimates of β ,
- is **bias-preventive** rather than corrective,
- is available in **Software** packages such as SAS, R, Stata...