



Can machine learning algorithms outperform traditional pricing methods?

Bor Harej

Bahnhofskolloquium, 6 January 2020



- Actuary SAA
- 12 years at Triglav, Head of QRM
- Joined Prime Re Solutions in 2019, shareholder

- President of SAA
- Board member of ASTIN
- Champion of ASTIN working party ICDML

- Purpose

Comparison of quality of fit of several algorithms calibrated on synthetic policy and claim data, where „true“ expected claims are known.

The presentation will focus on the results. The detailed explanation of individual algorithms is out of scope.

- Synthetic data

- Algorithms

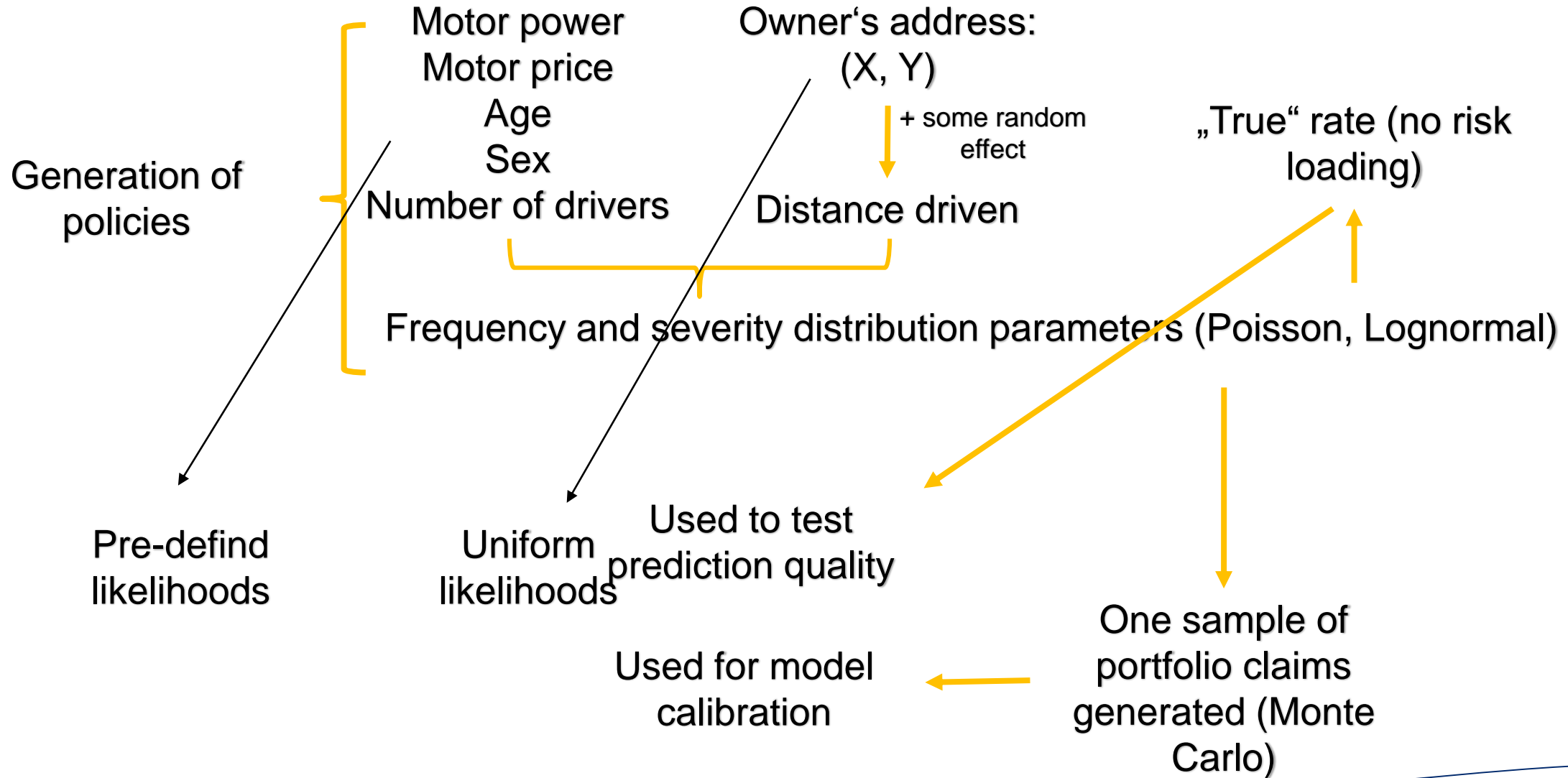
- Metrics

- **Results**

- Overfit check

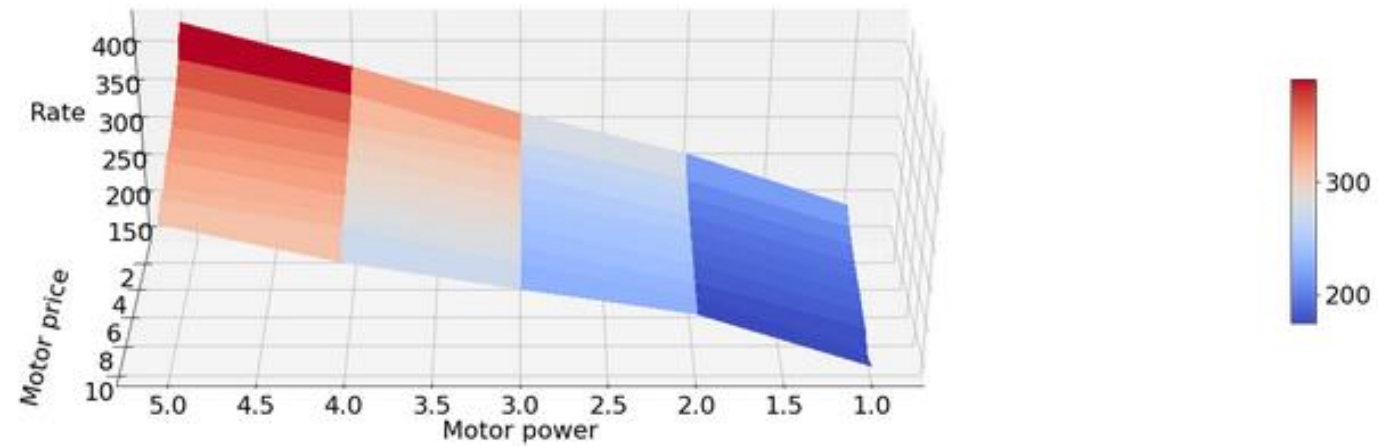
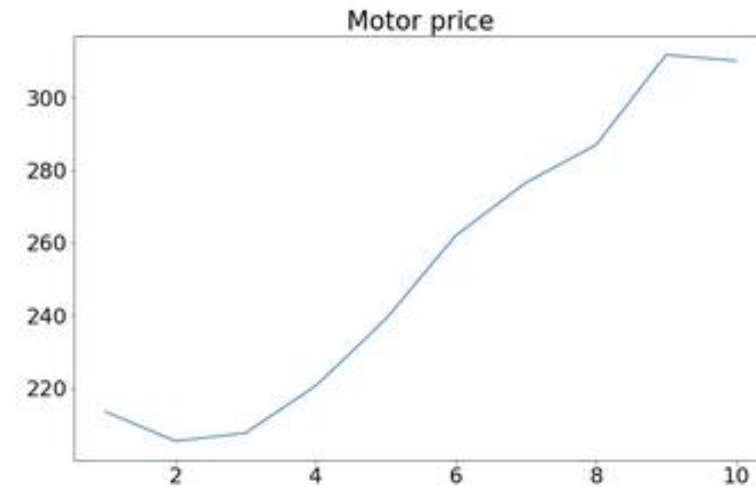
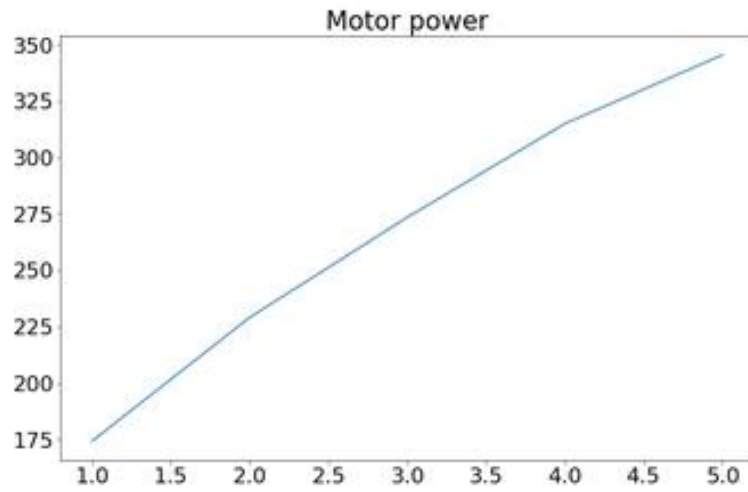


Synthetic data

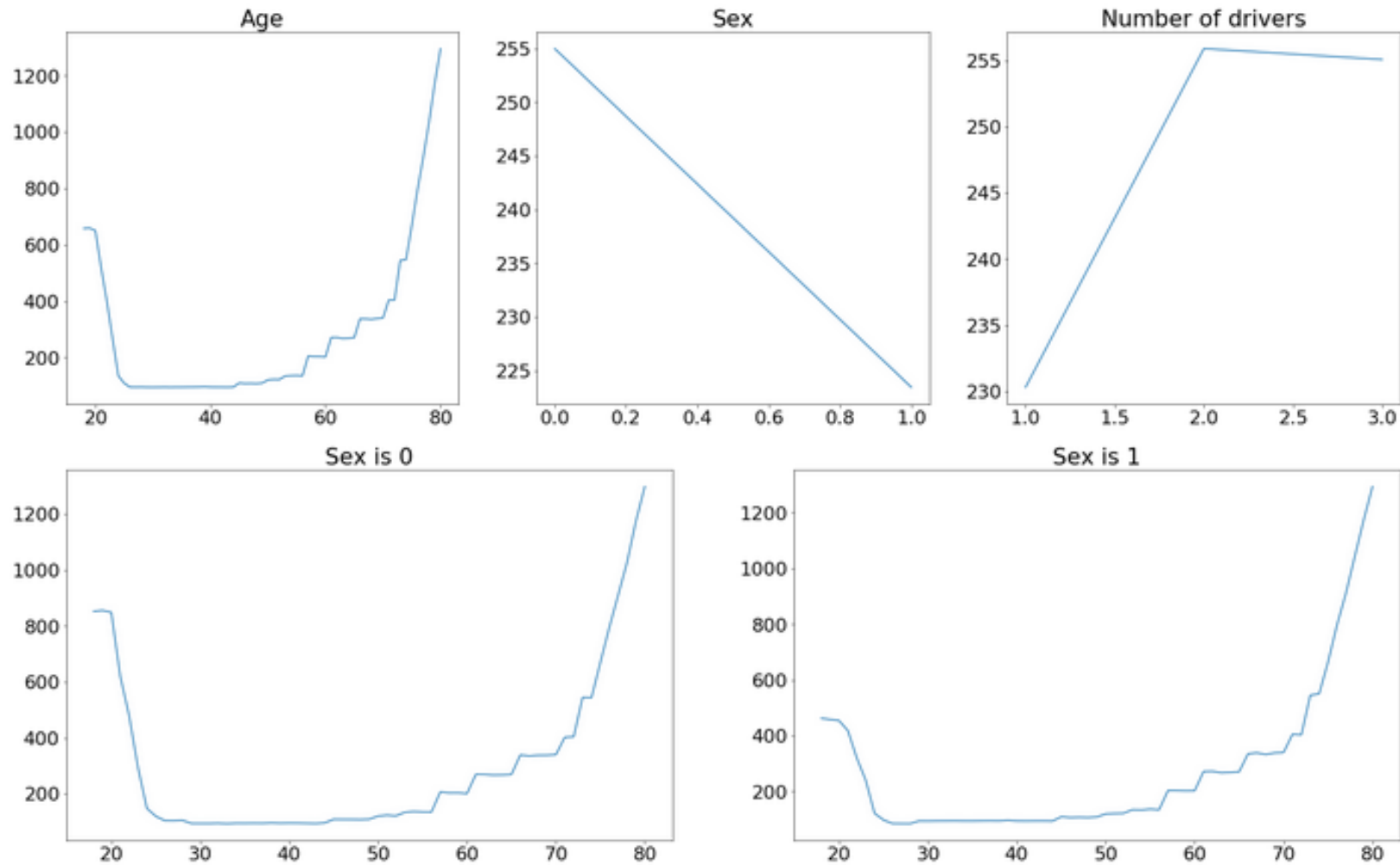


- Motor policy data (1 million)
 - Age (18-80), Sex (0,1), Motor power (1-5), Motor price (1-10), Number of drivers (1-3), owner's address (X, Y) (0-100,0-100)
 - Likelihoods defined or uniform
- Claims data
 - Poisson dist frequency, lognormal dist severity, total yearly claim used, lambda, mean and STD function of all 7 parameters
 - Function characteristics:
 - Non-linearity
 - Non-additivity
 - Non-linear dependency
 - Distance is a factor
 - Distance connected with address, some randomness added
- Simulated claims data + true policy premium rates (no risk premium)

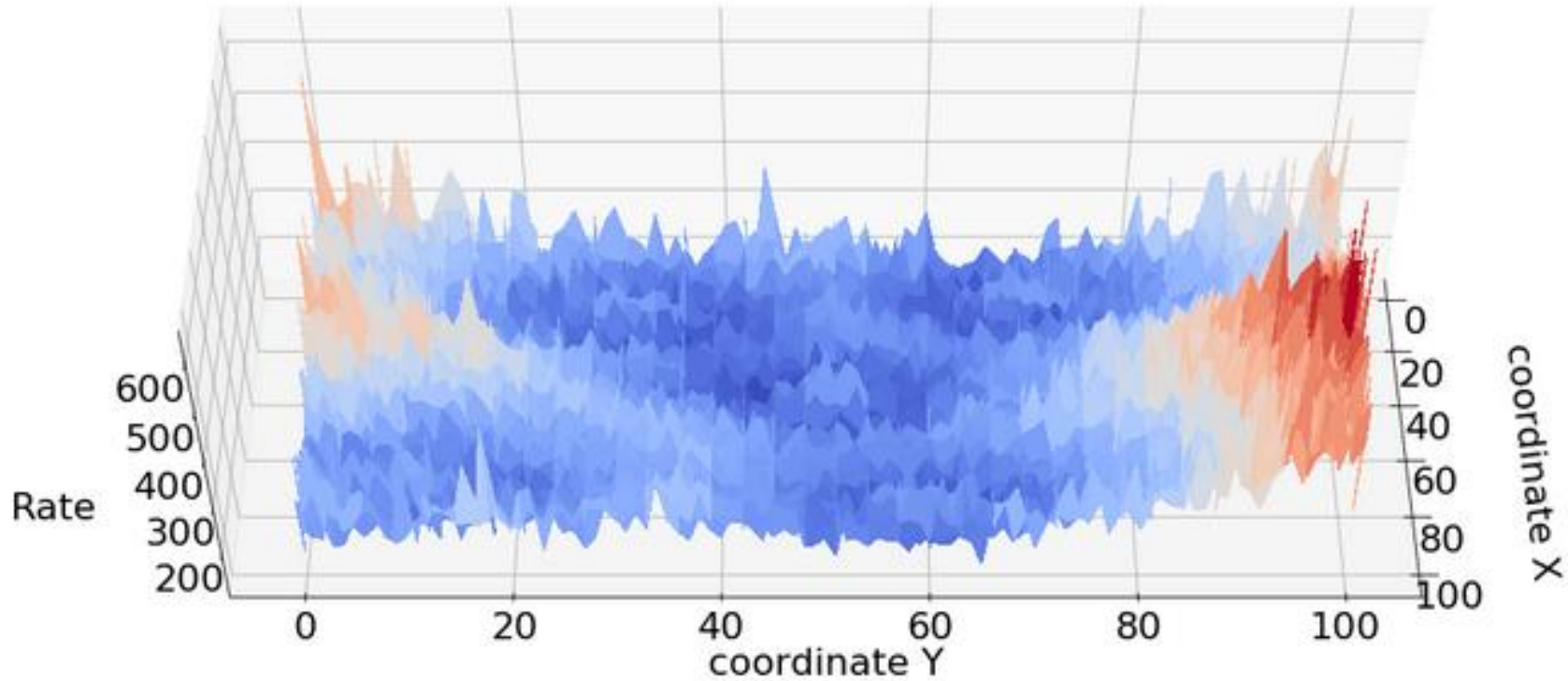
Synthetic data: dataset overview



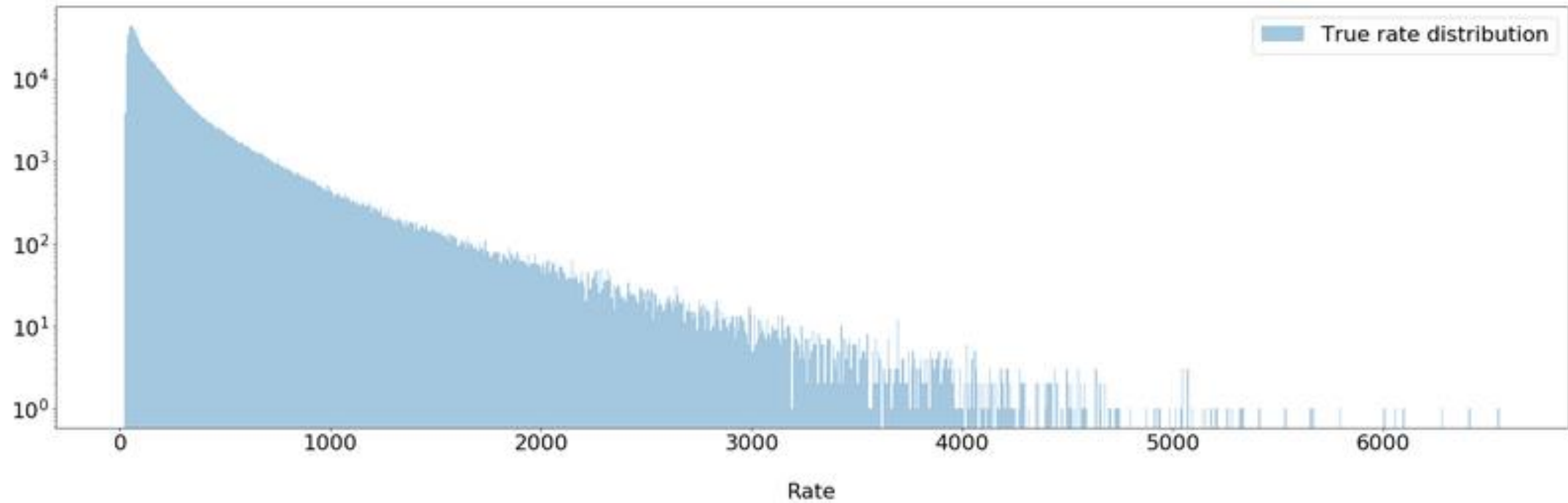
Synthetic data: dataset overview



Synthetic data: dataset overview



Synthetic data: rate distribution



- Traditionally used methods:
 - Generalized linear models (GLM)
 - Generalized additive models (GAM)

- Other machine learning algorithms:
 - Support vector machine (SVM)
 - Random forests (RF)
 - eXtreme gradient boosting (XGB)
 - Light gradient boosting (Light GBM)
 - Neural networks (NN)
 - Regression
 - Classification

kaggle™

- Charts
 - Averages by individual parameters
 - Averages by two parameters (3d)
 - Overall rate distribution

- Distance between „true“ rates and predicted rates
 - RMSE

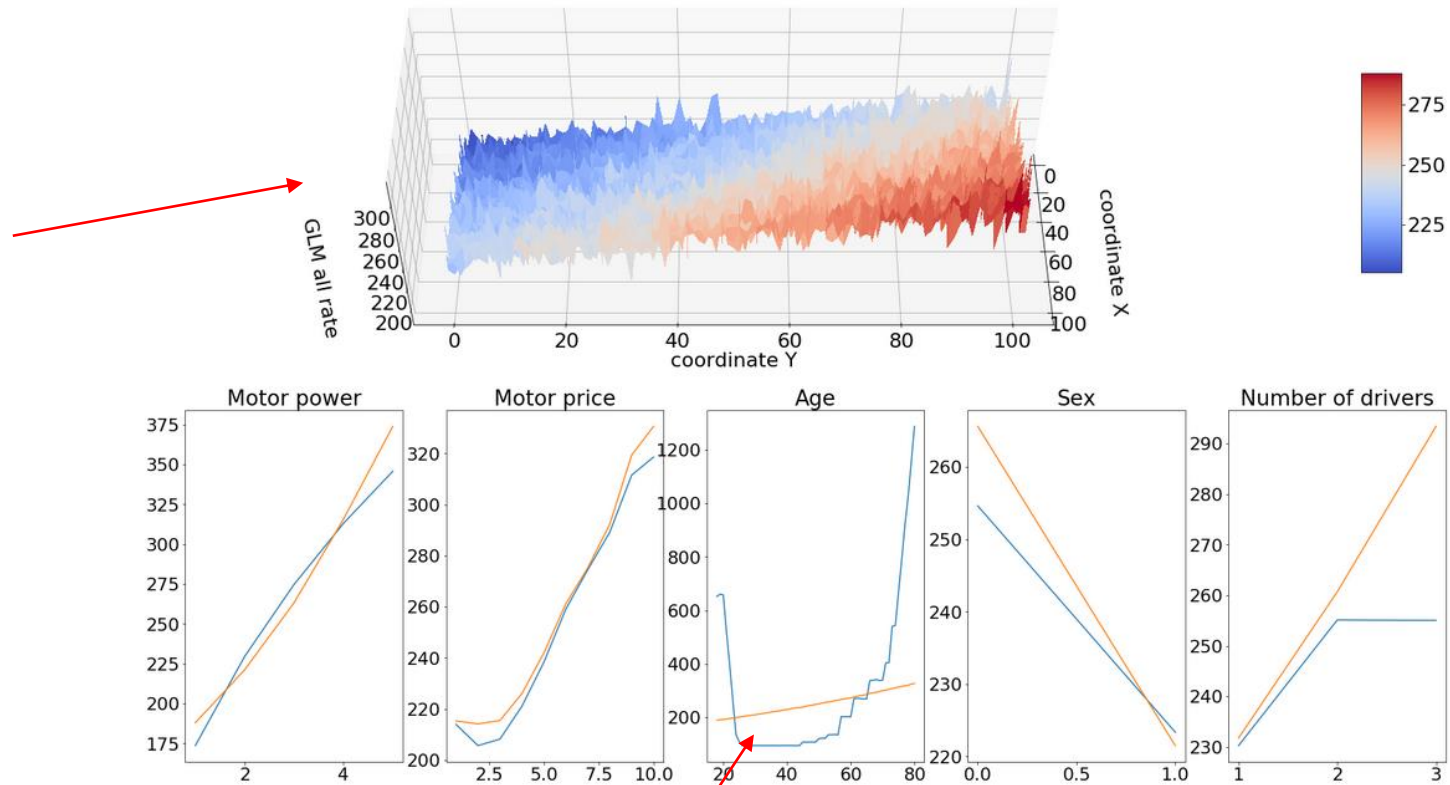
- Market share and profit
 - Assumptions:
 - Expected claim is a final premium rate
 - Cheapest option is taken

Traditionally used algorithms (GLM, GAM)

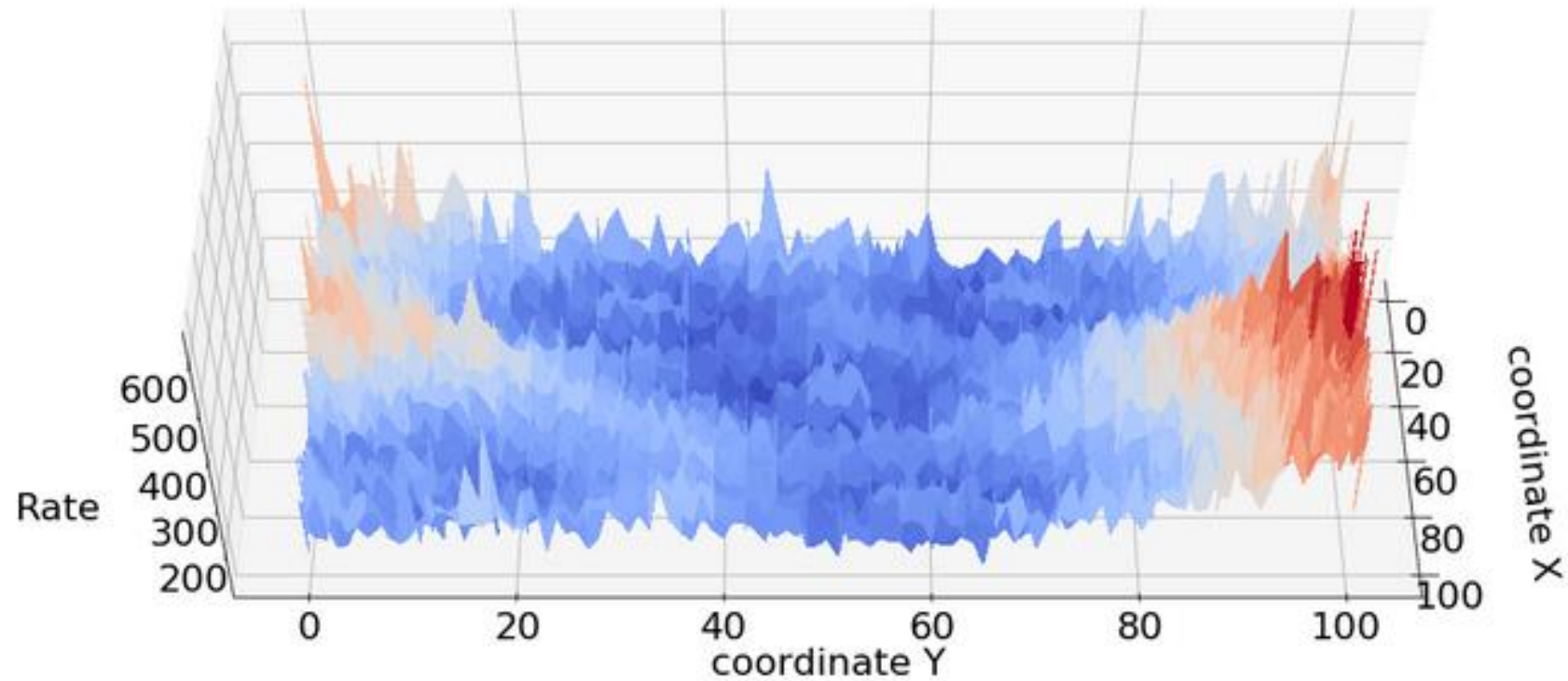
RESULTS

Results: GLM

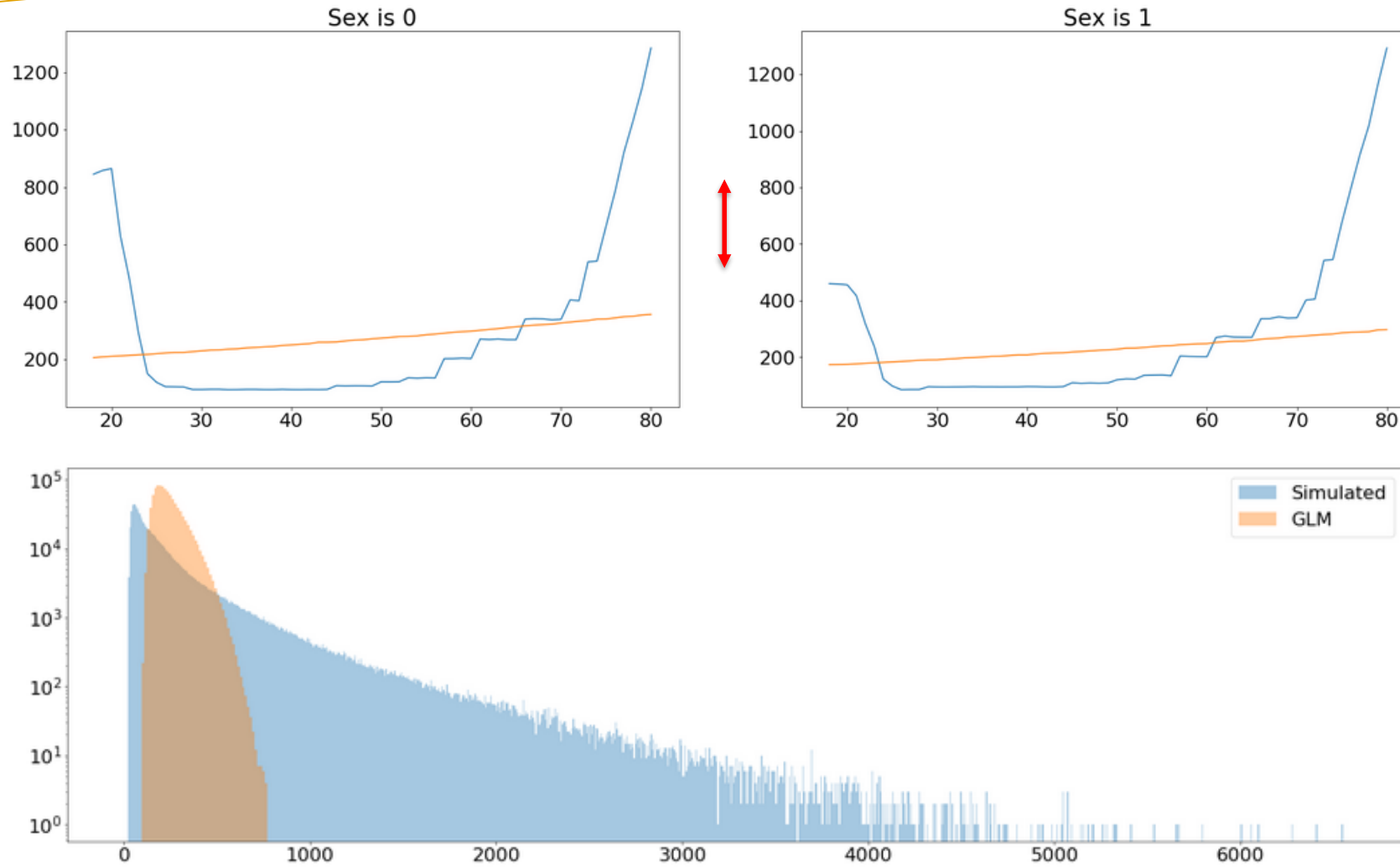
- Gamma distribution, log link, all parameters used



Synthetic data: premium factors

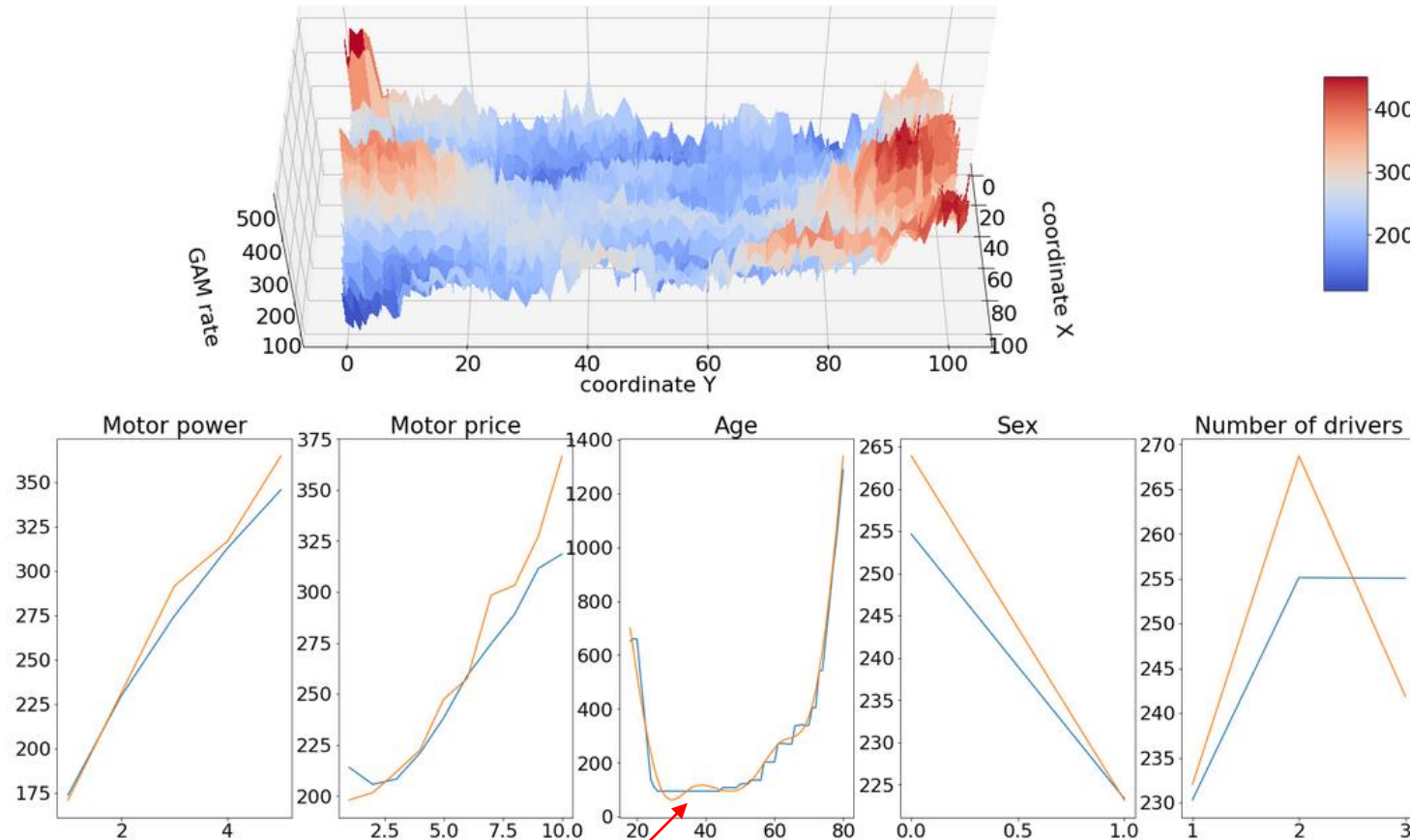


Results: GLM

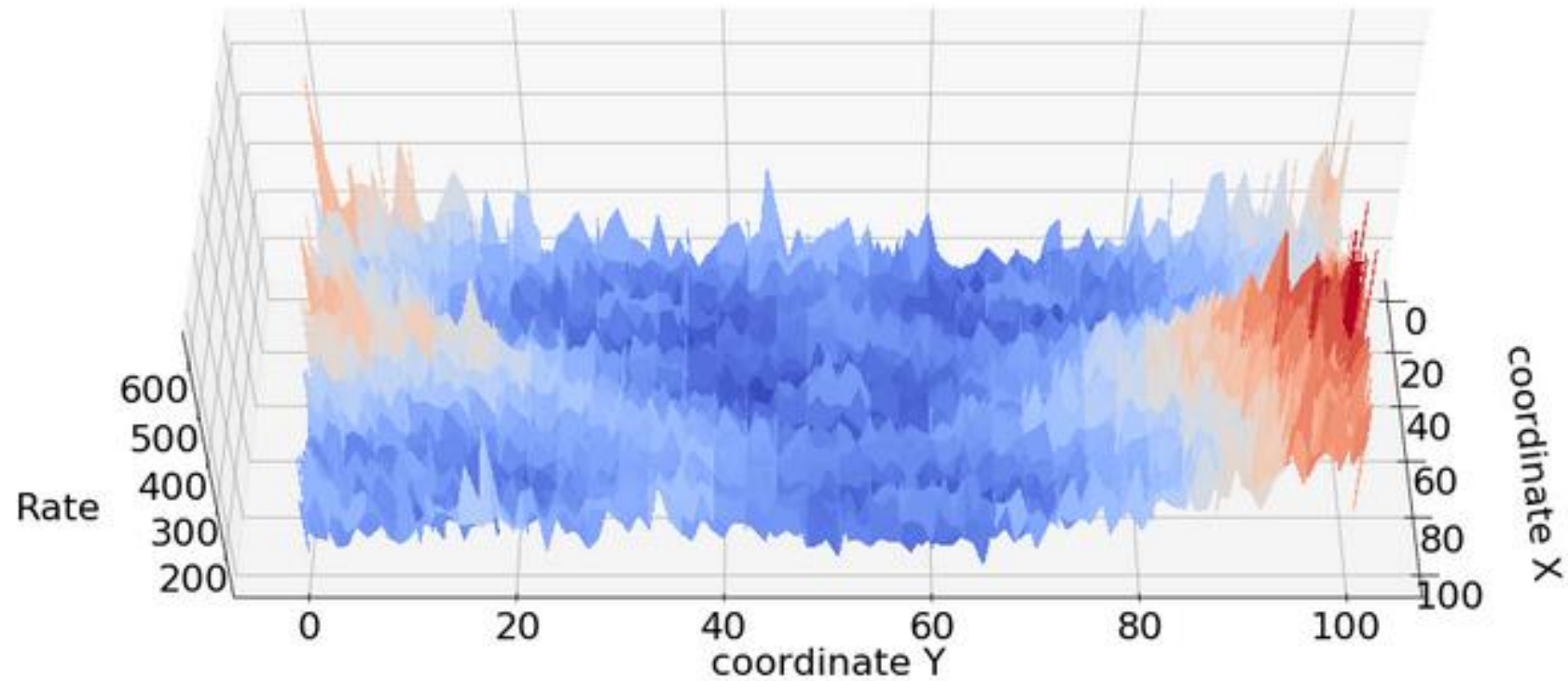


Results: GAM

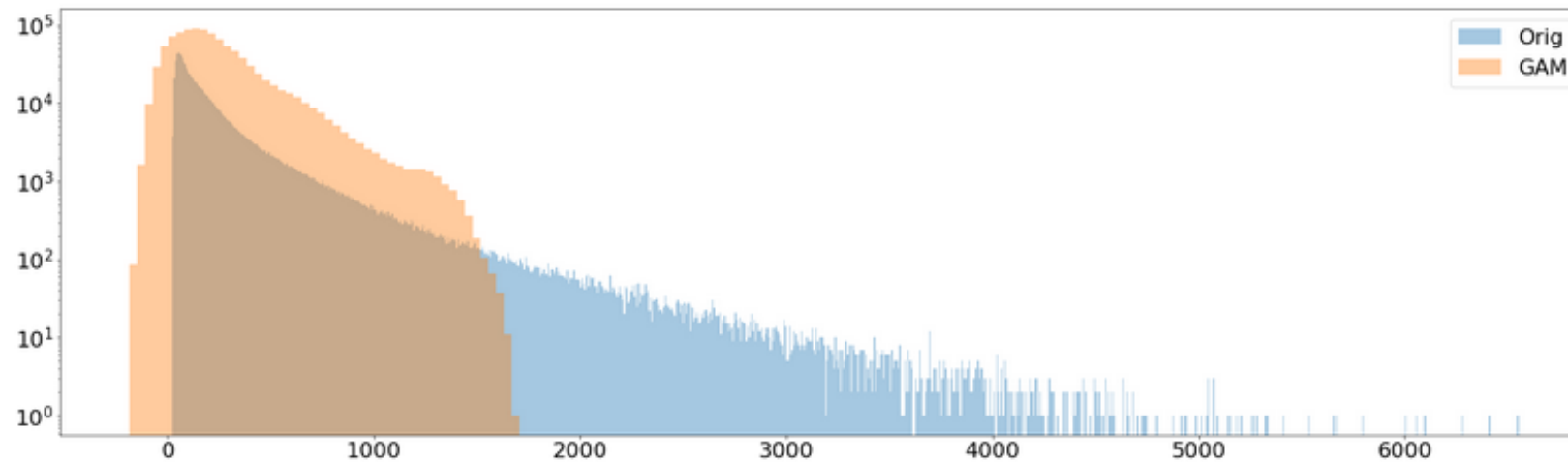
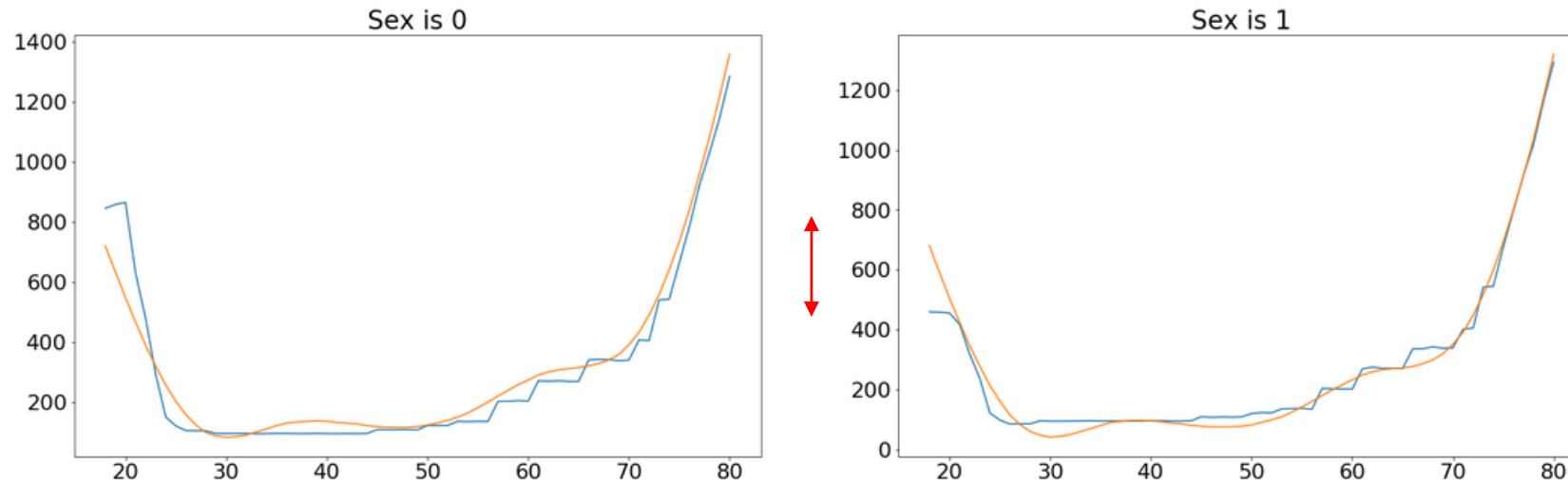
- 5 parameters with plain regression spline, coordinates with tensor product smooth, identity link function



Synthetic data: premium factors



Results: GAM

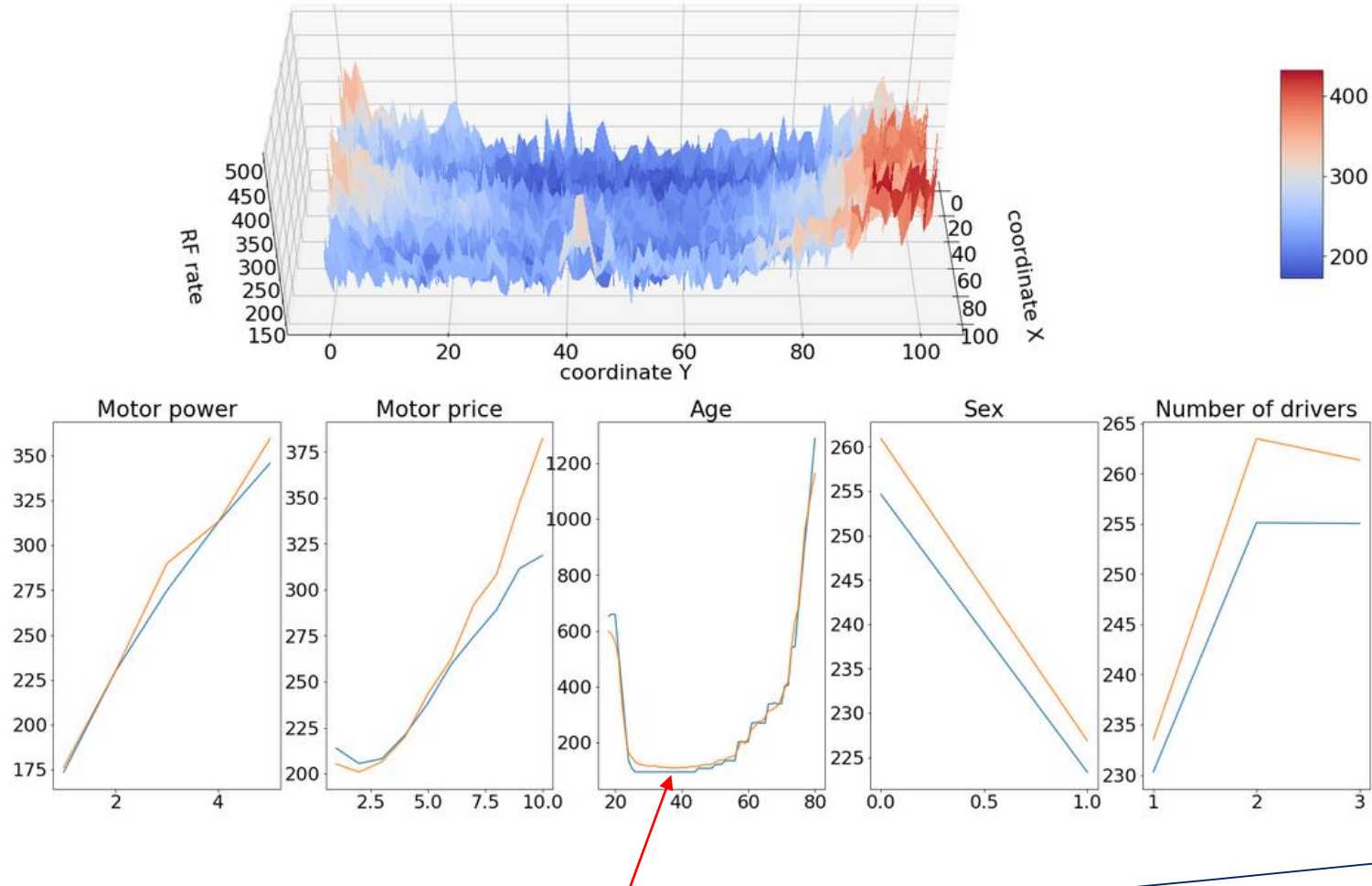


Other machine learning algorithms (RF, GB, NN)

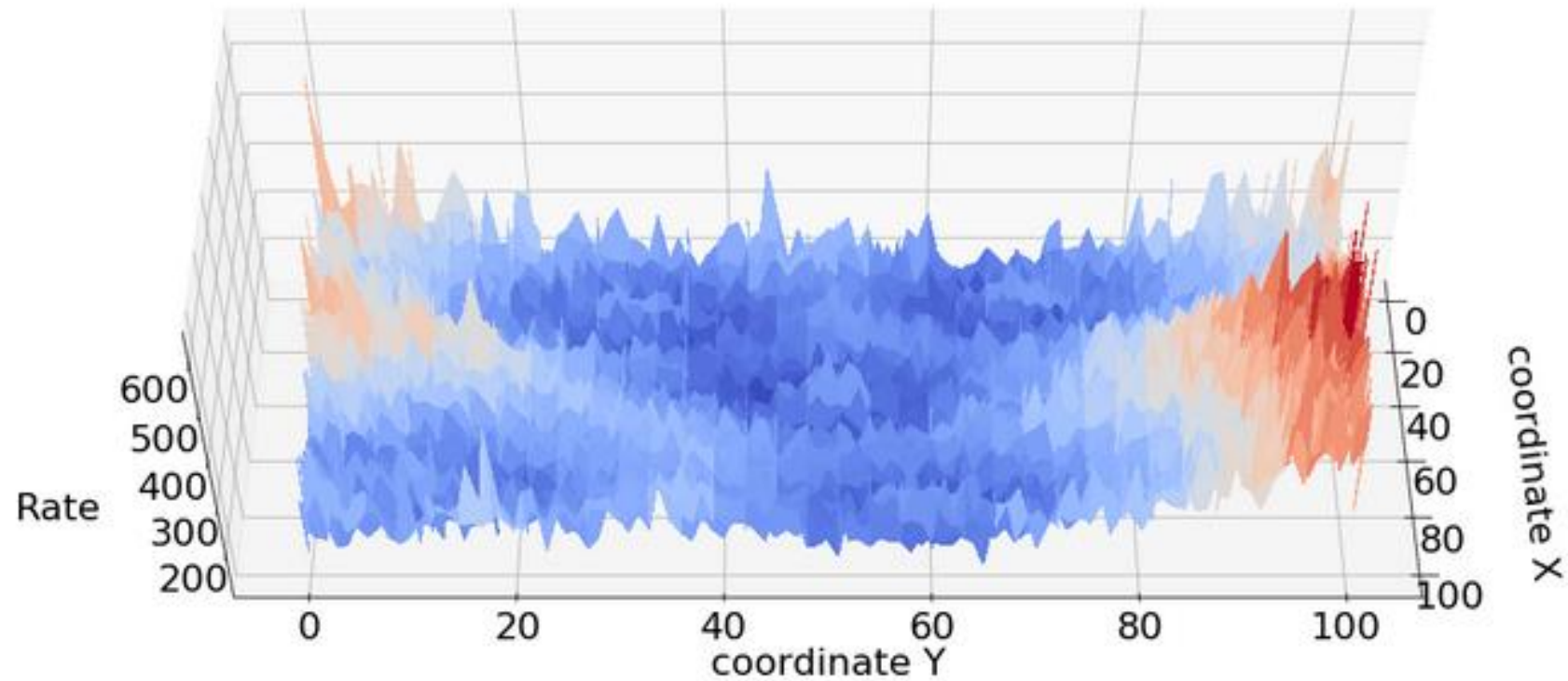
RESULTS

Results: Random forests

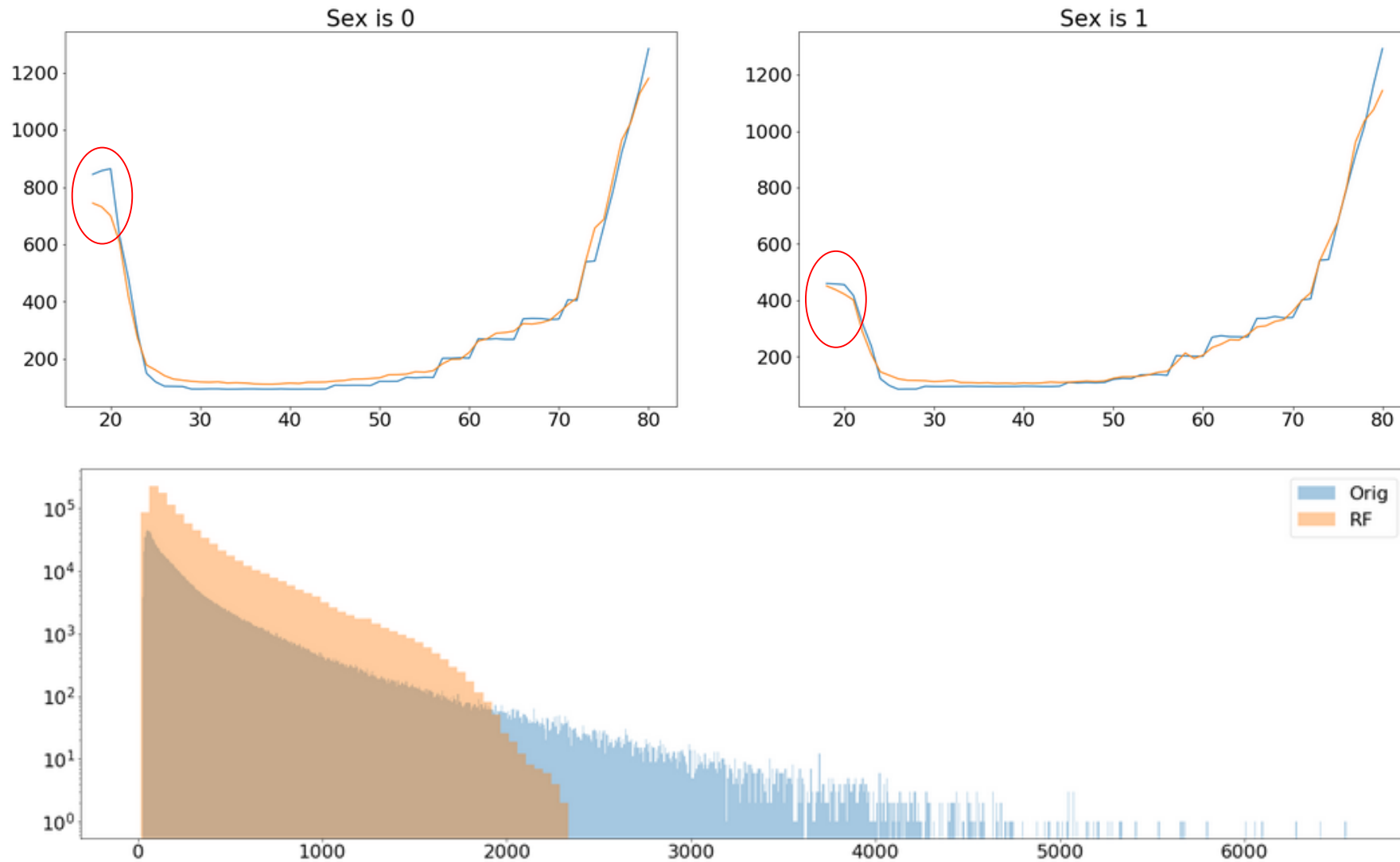
- Max depth, number of trees, min number of leafs, max features



Synthetic data: premium factors



Results: Random forests



XGB

- max depth, learning rate, no estimators, min child weight, gamma ...

Light GBM

- Max depth, learning rate, no estimators, boosting type, min data in leaf ...

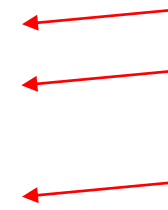
NN

- Architecture, activation functions, optimizer, loss function, dropout rate ...
- Classification classes

Results: overview



	(in 1.000 units)	Calibration
RMSE	GLM total predictions	296
	GAM total predictions	211
	Random forest predictions	202
	XGBoost predictions	195
	Light GBM predictions	194
	Neural networks predictions	205
	Classification neural networks predictions	193



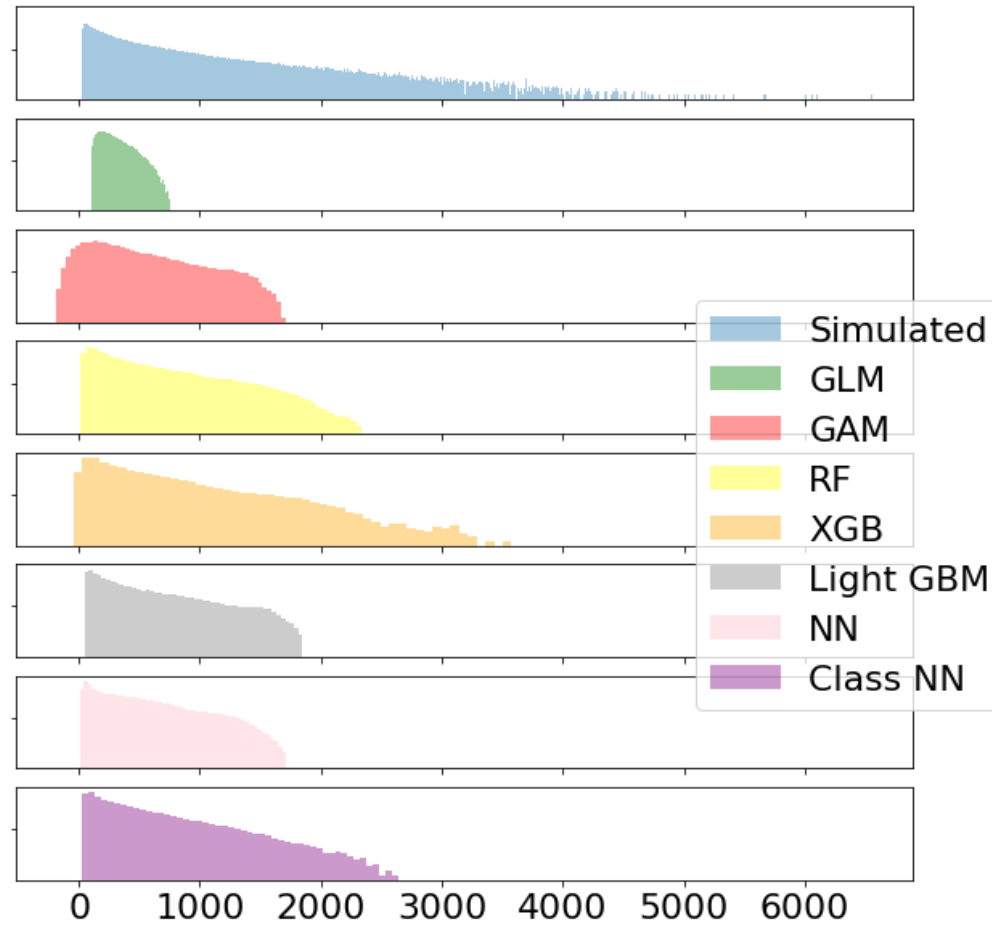
Qualitative assessment*:

	GLM	GAM	SVR	RF	XGB	Light GBM	NN	Class NN
Speed (learning)	4	2	1	3	3	4	2	2
Calibration complexity	4	3	4	3	2	3	2	2
Efficiency	2	3	1	4	5	5	4	5

* Subjective assessment, 1 worst, 5 best

Main issue: overfitting!

Results: rate distributions



Results: market share and profit



- Assumption: two players on market, one uses GLM for pricing, another classification NN

Total simulated claims	244.783.691	
Total (true) premium	238.995.840	-2,4%
Total premium GLM	241.323.027	-1,4%
Total premium Class NN	229.582.146	-6,6%

GLM vs Class NN	GLM premium	Share	True premium	Claims	Profit	Relative profit
Total premium	241.323.027	100%	238.995.840	244.783.691	-3.460.664	-1,4%
Premium below market price	65.323.200	27%	140.170.056	145.159.954	-79.836.754	-122,2%
Premium above market price	175.999.827	73%	98.825.783	99.623.738	76.376.089	43,4%
Class NN vs GLM	Class NN premium	Share	True premium	Claims	Profit	Relative profit
Total premium Class NN	229.582.146	100%	238.995.840	244.783.691	-15.201.545	-6,6%
Premium below market	95.056.387	41%	98.825.783	99.623.738	-4.567.351	-4,8%
Premium above market	134.525.759	59%	140.170.056	145.159.954	-10.634.195	-7,9%

Total market premium	160.379.587
Winner's curse ratio	67,1%

Results: market share and profit



Total claims	244.783.691
Total premium	238.995.840
Total premium GLM	241.323.027
Total premium RF	243.907.662
Total premium Light GBM	243.119.068
Total premium XGB	244.175.522
Total premium GAM	243.540.112
Total premium NN	237.157.861
Total premium Class NN	229.582.146

Relative profit ratio

party\market	GLM	GAM	RF	XGB	Light GBM	NN	Class NN
GLM		-86,4%	-101,3%	-105,7%	-108,1%	-94,1%	-122,2%
GAM	-6,9%		-32,7%	-31,0%	-34,9%	-21,1%	-43,8%
RF	1,5%	-9,0%		-12,3%	-15,7%	-10,3%	-21,3%
XGB	1,5%	-1,6%	-9,8%		-10,3%	-1,8%	-15,8%
Light GBM	2,7%	-3,7%	-3,6%	-4,5%		-3,1%	-11,8%
NN	-15,9%	-13,8%	-26,2%	-25,7%	-27,8%		-33,2%
Class NN	-4,8%	-8,9%	-13,2%	-12,6%	-14,5%	-6,8%	

Overfit check

- Re-run policy and claim simulations and use calibrated models for prediction

	(in 1.000 units)	Calibration	Validation
RMSE	GLM total predictions	296	296
	GAM total predictions	211	212
	Random forest predictions	202	203
	XGBoost predictions	195	196
	Light GBM predictions	194	195
	Neural networks predictions	205	206
	Classification neural networks	193	194

- **Traditionally pricing methods can be outperformed**
 - Example shown on synthetic data
- Different algorithms, might be tricky to calibrate and not to overfit
- Best fits: Light GBM, Classification NN
- Strong effect on profitability and market share
- Results seem to be stable

Bor Harej

bor.harej@prs-zug.com

QUESTIONS?