

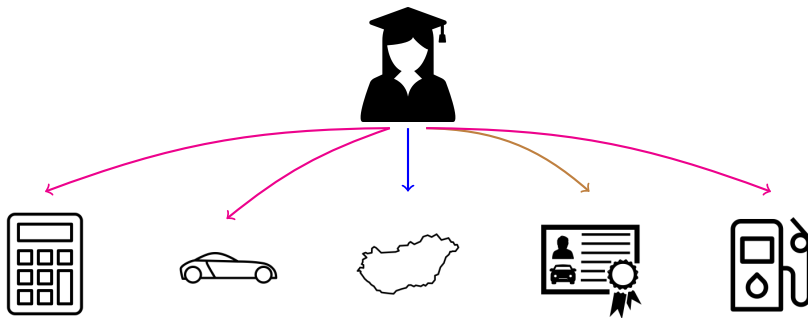
Sparse modeling of risk factors in insurance analytics

Sander Devriendt

Joint work with K. Antonio, T. Reynkens, E. Frees, R. Verbelen

eRum 2018, Budapest

May 15, 2018



Claim frequency and claim severity

as function of

nominal / numeric ~ ordinal / spatial

features

- ▶ Generalized Linear Models (GLMs) for frequency (\sim Poisson) and severity (\sim Gamma).
- ▶ How to:
 - (1) select risk factors or features?
 - (2) cluster (or bin or fuse) levels within a risk factor?
age groups / postal code clusters / clusters of car models
- ▶ Procedure should be data driven, scalable to large (big) data.
- ▶ End product is interpretable, within actuarial comfort zone.

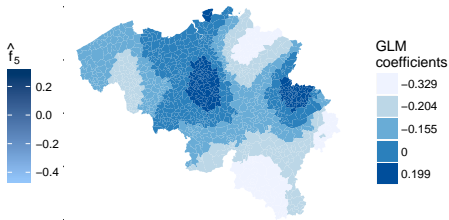
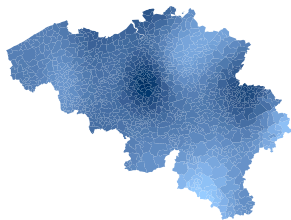
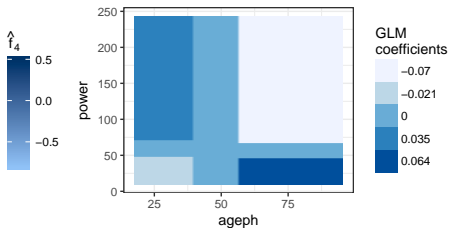
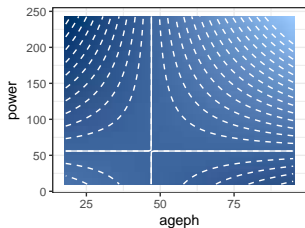
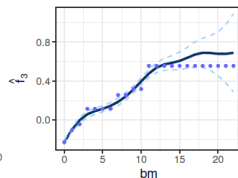
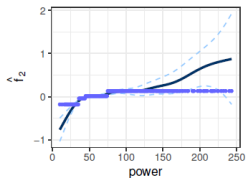
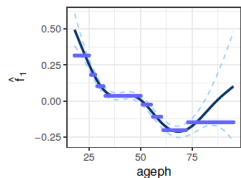
- ▶ Generalized Linear Models (GLMs) for frequency (\sim Poisson) and severity (\sim Gamma).
- ▶ How to:
 - (1) avoid overfitting with too many risk factors or levels?
 - (2) avoid underfitting with a priori binning/selection?

Henckaerts, Antonio et al., 2018 (accepted)

Stepwise procedure

- 1 Do an exhaustive search through variables to find best **GAM model**.
- 2 Use well-chosen **clustering algorithm** to bin 2D spatial effect.
- 3 Use **evolutionary trees** to bin 1D continuous effects and interactions.
- 4 **Fit GLM** with bins and clusters obtained in previous steps.

R packages: mgcv, classInt, evtree, rpart



Sparse modeling of risk factors in insurance analytics

Devriendt, Antonio, et al., 2018 (in progress)

LESS IS MORE

Ludwig Mies van der Rohe

► **Standard** GLM:

- **fit** data as good as possible,
- **no constraint** on parameters.



► **Regularized** GLM:

- **tradeoff** between fit and interpretability/sparsity/stability,
- **constraint** on parameters.

- **Less is more**: (Hastie, Tibshirani & Wainwright, 2015)

a sparse model is easier to estimate and interpret than a dense model.

- Regularize (with budget constraint t , or **regularization parameter λ**):

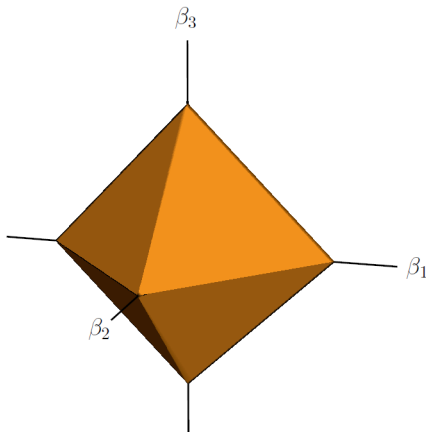
$$\min_{\beta_0, \beta} \{-\mathcal{L}(\beta_0, \beta)\} \text{ subject to } \|\beta\|_1 \leq t,$$

or equivalently

$$\min_{\beta_0, \beta} \left\{ -\mathcal{L}(\beta_0, \beta) + \lambda \cdot \sum_{j=1}^p |\beta_j| \right\}.$$

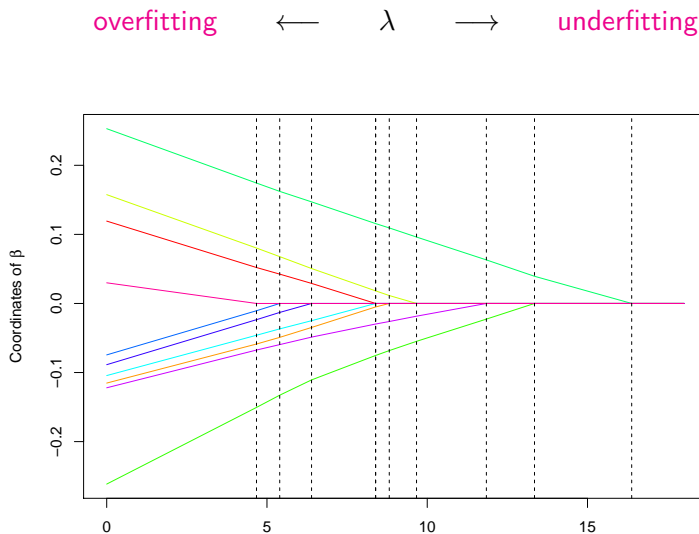
Shrinks coefficients and even sets some **to zero**.

Regularization = limited budget for $\beta_1, \beta_2, \beta_3$.



'Statistical Learning with Sparsity' - Hastie et al. (2015)

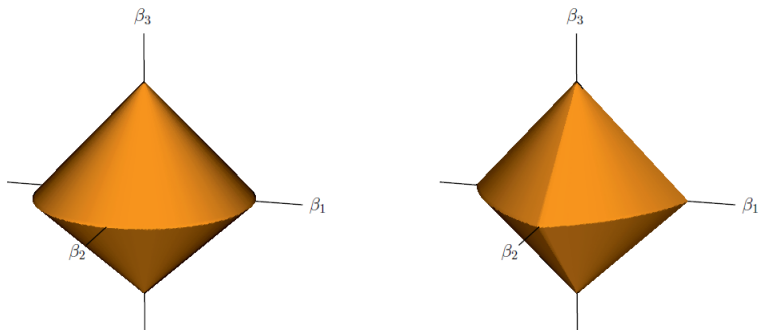
Package glmnet



- ▶ Adjust lasso regularization to the type of risk factor:
 - Determine type (nominal / numeric \sim ordinal / spatial);
 - Allocate logical penalty.
- ▶ Thus, for J risk factors, each with regularization term $P_j(\cdot)$, we want to optimize:

$$-\mathcal{L}(\beta_1, \dots, \beta_J) + \lambda \cdot \sum_{j=1}^J P_j(\beta_j).$$

Different variable type \rightarrow different penalty budget.

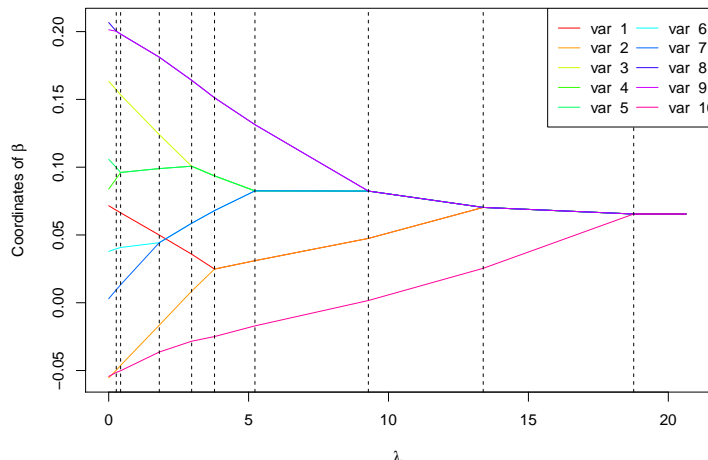


'Statistical Learning with Sparsity' - Hastie et al. (2015)

Package `genlasso`

overfitting $\leftarrow \lambda \rightarrow$ underfitting

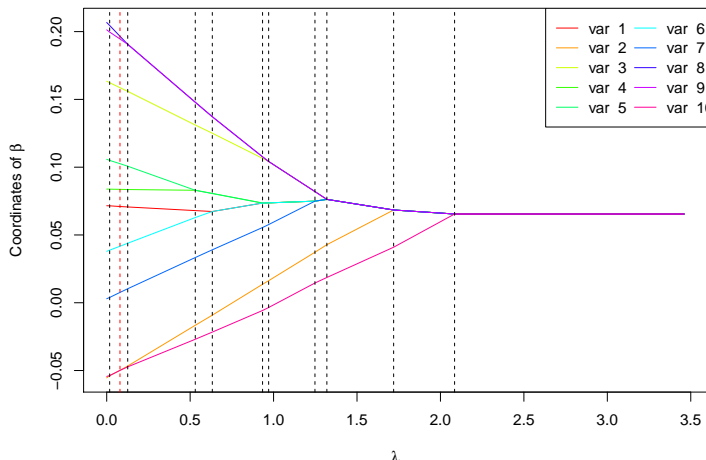
ordinal penalty example



Package genlasso

overfitting $\leftarrow \lambda \rightarrow$ underfitting

nominal penalty example



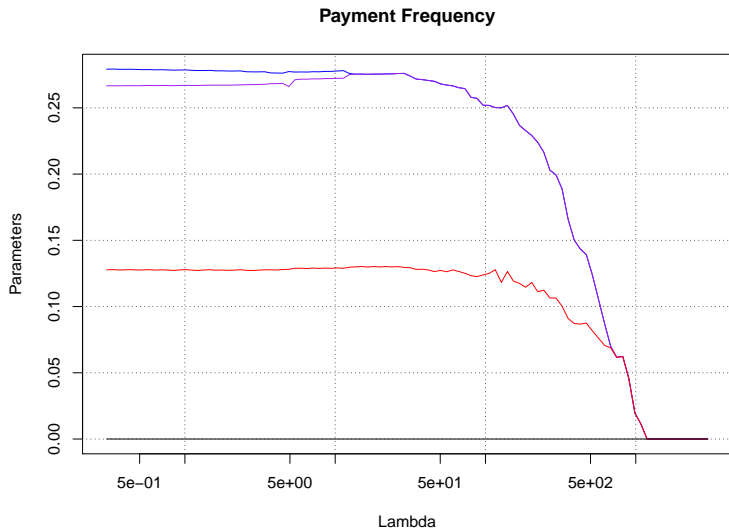
- ▶ Gertheiss & Tutz (2010) and Oelker & Gertheiss (2017):
 - GLMs with various penalties.
 - R package available: `gvcm.cat` (not maintained).
- ▶ Uses local quadratic approximations of penalties and PIRLS:
 - non-exact selection or fusion;
 - computationally intensive.

► Our contribution:

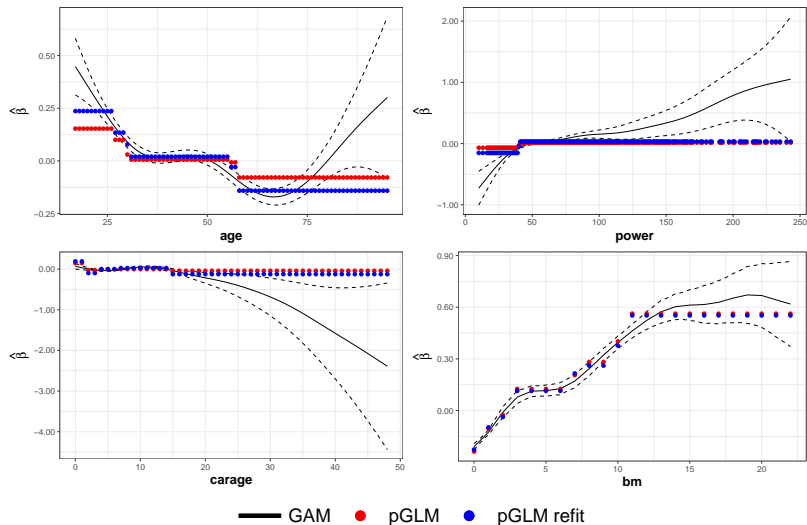
- implements an efficient algorithm (with proximal operators);
 - code bottleneck in C++ (Rcpp)
 - efficient linear algebra (RcppArmadillo)
 - parallel computations (parallel)
- scalable to big data (splits into smaller sub-problems);
- flexible regularization
 - penalty takes type of risk factor into account;
 - works for all popular penalties;

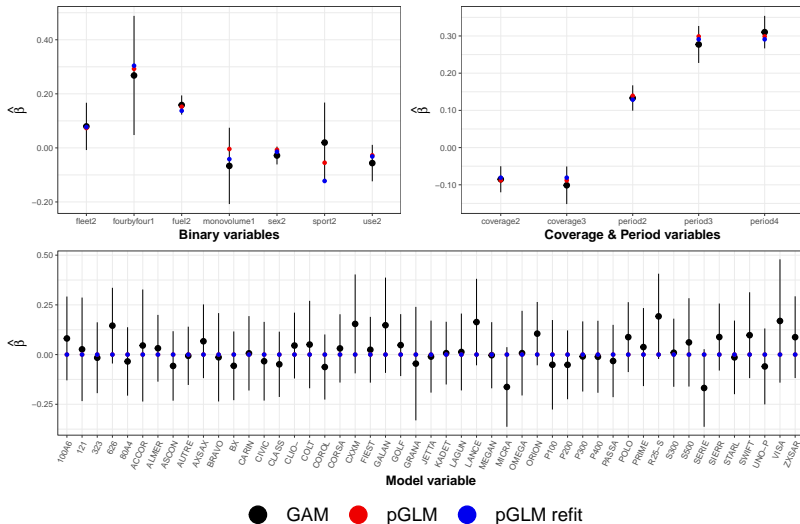
⇒ Package **mtppga** under construction.

- ▶ Frequency (and severity) information for $n = 163,234$ policyholders.
- ▶ 15 risk factors: binary, ordinal and nominal.
- ▶ Exposure modeled as offset.
- ▶ Fit Poisson GLM for frequency data with different penalties.



- ▶ Settings:
 - Incorporate **adaptive and standardization weights** for better consistency and predictive performance.
 - Tune λ with **out-of-sample MSE**
- ▶ **Re-estimate** the final sparse GLM with standard GLM routines (**from 146 to 30 params.**).





- ▶ less is more.
- ▶ Flexible regularization can help predictive modeling
- ▶ R package [mtp PGA](#) combines general framework with efficient algorithm.
- ▶ Package and working paper to be finalized.



- ▶ Tom Reynkens and colleagues
- ▶ You, the public

Henckaerts, R., Antonio, K., Clijsters, M. and Verbelen, R. (2018)
A data driven strategy for the construction of insurance tariff classes.
Scandinavian Actuarial Journal, published online.

Wood, S. (2006)
Generalized additive models: an introduction with R.
Chapman and Hall/CRC Press.

Gertheiss, J. and Tutz, G. (2010).
Sparse modeling of categorical explanatory variables.
The Annals of Applied Statistics, 4(4), 2150-2180.

Oelker, M. and Gertheiss, J. (2017).
A uniform framework for the combination of penalties in generalized
structured models.
Advances in Data Analysis and Classification, 11(1),97-120.

Parikh, N. and Boyd, S. (2013).

Proximal algorithms.

Foundations and Trends in Optimization, 1(3):123-231.

Hastie, T., Tibshirani, R. and Wainwright, M. (2015)

Statistical learning with sparsity: the Lasso and generalizations.

Chapman and Hall/CRC Press.