



Variable selection in general multinomial logit models



Gerhard Tutz^a, Wolfgang Pößnecker^{a,*}, Lorenz Uhlmann^b

^a Department of Statistics, Ludwig-Maximilians-University, 80539 München, Germany

^b Department of Medical Biometry, Ruprecht-Karls-University, 69120 Heidelberg, Germany

ARTICLE INFO

Article history:

Received 21 October 2013

Received in revised form 8 September 2014

Accepted 10 September 2014

Available online 21 September 2014

Keywords:

Logistic regression

Multinomial logit model

Variable selection

Lasso

Group Lasso

CATS Lasso

ABSTRACT

The use of the multinomial logit model is typically restricted to applications with few predictors, because in high-dimensional settings maximum likelihood estimates tend to deteriorate. A sparsity-inducing penalty is proposed that accounts for the special structure of multinomial models by penalizing the parameters that are linked to one variable in a grouped way. It is devised to handle general multinomial logit models with a combination of global predictors and those that are specific to the response categories. A proximal gradient algorithm is used that efficiently computes stable estimates. Adaptive weights and a refitting procedure are incorporated to improve variable selection and predictive performance. The effectiveness of the proposed method is demonstrated by simulation studies and an application to the modeling of party choice of voters in Germany.

© 2014 Elsevier B.V. All rights reserved.

1. Introduction

The multinomial logit model is the most frequently used model in regression analysis for un-ordered multi-category responses. The maximum likelihood (ML) method, which is typically used for estimation, has the drawback that it requires more observations than parameters to be estimated. The amount of parameters, however, increases rapidly when the number of predictors grows, as multinomial logit models employ several coefficients for each explanatory variable. Therefore, ML estimates tend to deteriorate quickly and interpretability suffers as well, so that the number of predictors in the model is severely limited. For these reasons, variable selection is necessary to obtain multinomial logit models that are both interpretable and reliable.

As a motivating example, our application uses data from the German Longitudinal Election Study (GLES) about party choice of voters during the 2009 parliamentary elections for the German Bundestag. Modeling the decision of voters for specific political parties and determining the major factors behind their preference are of great interest in political sciences. The available parties to choose from are the Christian Democratic Union (CDU), the Social Democratic Party (SPD), the Green Party (Bündnis 90/Die Grünen), the Liberal Party (FDP) and the Left Party (Die Linke). As explanatory variables, various individual characteristics of the voter are considered, like, for example, gender, age or education, see Section 5 for a complete list. The main goal of our analysis is to select those predictors that influence party choice and to remove the rest. Besides improving interpretability, this is beneficial for polling firms and in opinion research. If, for example, gender was found to be irrelevant for party preference, one could save time and money while performing opinion polls as one would not have to care about a representative gender ratio among the interviewed persons.

* Correspondence to: Seminar of Applied Stochastics, Department of Statistics, Ludwig-Maximilians-University Munich, Ludwigstraße 33, 80539 München, Germany.

E-mail address: wolfgang.poessnecker@stat.uni-muenchen.de (W. Pößnecker).

While variable selection for such individual-specific predictors requires specific methodology due to the particular structure of multinomial logit models, this dataset offers another challenge in the form of predictors that are party-specific. For various topics like immigration or nuclear energy, participants of the study were asked how they perceive the parties' stance on this issue. Additionally, they stated their personal position on the topic. From this information, the distance between the personal point of view and the perceived position of the parties can be computed. These distances are then included into the model. Since they take different values for different categories of the response, they are an example of so-called category-specific predictors. To be able to deal with the challenges of such a dataset, a method for variable selection in such general multinomial logit models is developed in this paper. It accounts for the categorical and multivariate structure of these models and works with both global and category-specific predictors.

The standard method for variable selection are forward/backward strategies which have been used for a long time, but are notoriously unstable and computationally costly, so that they cannot be recommended (see, for example, Hastie et al., 2009). An established alternative is penalty approaches for regularized variable selection. For linear and generalized linear models (GLMs), a variety of such methods has been proposed. The most prominent example is the Lasso (Tibshirani, 1996) and its extensions to Fused Lasso (Tibshirani et al., 2005) and Group Lasso (Yuan and Lin, 2006). Alternative regularized estimators that enforce variable selection are the Elastic Net (Zou and Hastie, 2005), SCAD (Fan and Li, 2001), the Dantzig selector (Candes and Tao, 2007) and boosting approaches (Bühlmann and Yu, 2003; Bühlmann and Hothorn, 2007; Tutz and Binder, 2006).

These methods, however, were developed for models with univariate response. Because the multinomial logit model is not a common univariate GLM, these methods cannot be applied directly. As mentioned previously, the effect of one predictor variable is represented by several parameters. Therefore, one has to distinguish between variable selection and parameter selection, where variable selection is only achieved if all effects/parameters that belong to one variable are simultaneously removed from the model. Early suggestions for regularization in multinomial logit models (Krishnapuram et al., 2005; Friedman et al., 2010) use L_1 -type penalties that shrink all the parameters individually. Thus, they do not use the natural grouping of coefficients that is available, with each group containing the parameters that belong to the same explanatory variable. In particular, they pursue the goal of parameter selection and cannot directly promote variable selection.

Therefore, a more explicit approach to variable selection in multinomial logit models is to use a grouped penalization with groups that consist of all coefficients belonging to the same predictor variable. Such a grouped variable selection was used in general multivariate regression, among others, by Turlach et al. (2005) and Argyriou et al. (2007). The present paper extends and complements more recent work by Simon et al. (2013), Vincent and Hansen (2014) and Chen and Li (2013), who were among the first to explicitly embed the idea of the group lasso into the multinomial regression framework. However, all three of these papers only consider global predictors that are constant/independent from the observed response category. By contrast, we consider the general multinomial logit model in which a mix of global and category-specific predictors is used. Moreover, we extend our penalization approach to the case of categorical predictors with more than two categories. Since all dummy variables belonging to such a predictor should be selected jointly, this creates an additional level of grouping on top of the grouping across response categories. Additionally, the concept of adaptive weights (Zou, 2006; Wang and Leng, 2008) is adopted. We demonstrate that this distinctly improves variable selection and prediction accuracy. From an algorithmic point of view, the three aforementioned papers all use variations of a coordinate descent algorithm for the computation of numerical estimates. By contrast, we use the fast iterative shrinkage thresholding algorithm (FISTA) of Beck and Teboulle (2009), which is an accelerated version of proximal gradient methods.

This paper is organized as follows: In Section 2 we introduce the general multinomial logit model and suggest a penalty that yields proper variable selection. Extensions to incorporate category-specific variables and multi-categorical predictors both in the model and in the penalization are discussed separately. Regularized estimation is considered in Section 3, where a proximal gradient algorithm is derived that efficiently solves the corresponding estimation problem. The performance of our estimator is investigated in simulation studies in Section 4. Then, in Section 5, the real data example from the German Longitudinal Election Study is analyzed using the developed methodology.

2. Model and regularization

2.1. The multinomial logit model with category-specific covariates

For data (y_i, \mathbf{x}_i) , $i = 1, \dots, n$, with y_i denoting an observation of the categorical response variable $Y \in \{1, \dots, k\}$ and \mathbf{x}_i the p -dimensional vector of predictors, the multinomial logit model in its generic form specifies

$$\pi_{ir} = P(Y = r | \mathbf{x}_i) = \frac{\exp(\beta_{r0} + \mathbf{x}_i^T \boldsymbol{\beta}_r)}{\sum_{s=1}^k \exp(\beta_{s0} + \mathbf{x}_i^T \boldsymbol{\beta}_s)} = \frac{\exp(\eta_{ir})}{\sum_{s=1}^k \exp(\eta_{is})}, \quad (1)$$

where $\boldsymbol{\beta}_r^T = (\beta_{r1}, \dots, \beta_{rp})$. Since parameters $\beta_{10}, \dots, \beta_{k0}, \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_k$ are not identifiable, additional constraints are needed. Typically, one of the response categories is chosen as reference category. We use category k as the reference category by

setting $\beta_{k0} = 0$, $\beta_k = \mathbf{0}$. With this choice, the linear predictors η_{ir} , $r = 1, \dots, k-1$, correspond to the log odds between category r and the reference category k .

The model given in (1) is the most commonly used form of the multinomial logit model. However, it only uses “global” predictors that do not vary over response categories, a restriction that is not always appropriate in practice. In particular in the modeling of choice, when an individual chooses among alternatives $1, \dots, k$, one wants to model the effects of characteristics of the individual like age and gender, but also account for measured attributes of the alternatives $1, \dots, k$. In the modeling of preference for parties the alternatives are characterized by positions on policy dimensions. When the choice refers to transportation mode, the potential attributes are price and duration, which vary across the alternatives and therefore are category-specific. Then, in addition to the global predictors \mathbf{x}_i , a set of category-specific predictors $\mathbf{w}_{i1}, \dots, \mathbf{w}_{ik}$ is available, where \mathbf{w}_{ir} contains the attributes of category r . Including category-specific predictors \mathbf{w}_{ir} into the model yields the general multinomial logit model.

With k as the reference category, the set of linear predictors generalizes to

$$\eta_{ir} = \beta_{r0} + \mathbf{x}_i^T \beta_r + (\mathbf{w}_{ir} - \mathbf{w}_{ik})^T \alpha, \quad r = 1, \dots, k-1. \quad (2)$$

The second term specifies the effect of the global variables, and the third term specifies the effect of the difference $\mathbf{w}_{ir} - \mathbf{w}_{ik}$ on the choice between category r and the reference category. For the interpretation of parameters it is often useful to consider this choice the consequence of underlying latent utilities associated with the different alternatives. Let the latent utility of person i be specified by $u_{ir} = \gamma_{r0} + \mathbf{x}_i^T \gamma_r + \mathbf{w}_{ir}^T \alpha$ and the corresponding random utility be given as $U_{ir} = u_{ir} + \varepsilon_{ir}$, where ε_{ir} follows the Gumbel distribution, which has distribution function $F(\varepsilon) = \exp(-\exp(-\varepsilon))$. It is assumed that each individual chooses the alternative that yields maximum utility, that is, the link between the observable choice Y_i and the unobservable random utility is given by $Y_i = r \Leftrightarrow U_{ir} = \max_{j=1, \dots, k} U_{jr}$. For the probabilities of the alternatives, one then obtains the multinomial logit model with predictor

$$\eta_{ir} = u_{ir} - u_{ik} = (\gamma_{r0} - \gamma_{k0}) + \mathbf{x}_i^T (\gamma_r - \gamma_k) + (\mathbf{w}_{ir} - \mathbf{w}_{ik})^T \alpha,$$

which has the form given in (2) with $\beta_{r0} = (\gamma_{r0} - \gamma_{k0})$, $\beta_r = (\gamma_r - \gamma_k)$ (see [Yellott, 1977](#); [McFadden, 1973](#)). An extensive discussion of the multinomial logit model as a multivariate GLM is given in [Tutz \(2012\)](#).

2.2. Regularization: the Categorically Structured Lasso

We first focus on regularization in multinomial logit models with only global predictors in which the overall parameter vector is given by $\beta^T = (\beta_{10}, \dots, \beta_{k-1,0}, \beta_1^T, \dots, \beta_{k-1}^T)$. A common way to regularize a model is penalty approaches in which one maximizes the penalized log-likelihood

$$l_p(\beta) = l(\beta) - \lambda J(\beta),$$

where $l(\beta)$ denotes the usual log-likelihood, λ is a tuning parameter and $J(\beta)$ is a functional that typically penalizes the size of the parameters. While the tuning parameter determines the strength of the regularization, the functional determines the properties of the penalized estimator. [Tibshirani \(1996\)](#) introduced the Lasso that penalizes the L_1 -norm of coefficients, that is, $J(\beta) = \|\beta\|_1$, and showed that it facilitates sparse solutions in which coefficients of weak predictors are set to exactly zero. For the multinomial logit model without category-specific predictors, direct application of the Lasso yields the L_1 -norm of β :

$$J(\beta) = \|\beta\|_1 = \sum_{r=1}^{k-1} \|\beta_r\|_1 = \sum_{r=1}^{k-1} \sum_{j=1}^p |\beta_{rj}|. \quad (3)$$

[Friedman et al. \(2010\)](#) used the slightly more general Elastic Net penalty, which is an adaptation of the original Elastic Net ([Zou and Hastie, 2005](#)) to multinomial response models. In terms of variable selection, Elastic Net and Lasso both share the same drawback that selection focuses on parameters but not on variables. In univariate regression models, for which these penalties were developed, this distinction is irrelevant if only continuous or binary predictors are used because then each predictor influences the response through only one coefficient. By contrast, multinomial logit models use a whole vector $\beta_{\bullet j} = (\beta_{1j}, \dots, \beta_{k-1,j})$ of parameters to capture the influence of predictor x_j . The ordinary Lasso penalty from Eq. (3), however, only encourages selection of single parameters β_{rj} . Thus, variable selection is only achieved if $k-1$ coefficients are simultaneously shrunk to zero, but the Lasso does not enforce such behavior.

Therefore, the ordinary Lasso is not ideal for variable selection in multinomial models. Although one might suspect that setting many coefficients to zero still improves interpretability, it is seen from Eq. (1) that the probability of all response categories is influenced by a predictor x_j if just one of the corresponding coefficients β_{rj} , $r = 1, \dots, k-1$ is non-zero. Hence, there is a strong incentive in multinomial models to perform true variable selection by simultaneously removing all effects of a predictor from the model.

Therefore, we suggest to penalize the groups $\beta_{\bullet j}$ of parameters that are linked to one variable. For simplicity, we assume all predictors to be metric, standardized and centered around zero. Extensions to categorical predictors follow in Section 2.3.

If no category-specific predictors are included, we will use the penalty

$$J(\beta) = \sum_{j=1}^p \|\beta_{\bullet j}\| = \sum_{j=1}^p (\beta_{1j}^2 + \cdots + \beta_{k-1,j}^2)^{1/2}, \quad (4)$$

where $\|u\| = \|u\|_2 = \sqrt{u^T u}$ denotes the L_2 -norm. This penalty enforces variable selection, that is, all the parameters in $\beta_{\bullet j}$ are simultaneously shrunk toward zero. It is strongly related to the Group Lasso (Yuan and Lin, 2006; Meier et al., 2008). However, in the Group Lasso the grouping refers to the parameters that are linked to a categorical predictor within a univariate regression model whereas in the present model grouping arises from the multivariate structure of the multinomial logit model. An approach similar to penalty (4) was recently advocated in Simon et al. (2013), Vincent and Hansen (2014) and Chen and Li (2013).

Besides the difference between unstructured selection of parameters and structured selection of variables, there is another advantage of the grouped penalty over the ordinary Lasso: Although neither penalty (3) nor penalty (4) are invariant to the choice of a reference category, in our experience the sparsity induced by the Lasso is strongly influenced by this choice while the sparsity of the grouped approach is typically not affected. An empirical illustration of this trend is given for the party choice data in Table 2, Section 5. In this application the grouped penalty produces stable sparsity while the number of nonzero coefficients for Lasso solutions differs by up to 29% when using alternative reference categories. This instability of the ordinary Lasso is also the main reason why we chose not to combine the ordinary lasso with our grouped approach in a “sparse group lasso”-like fashion as in Vincent and Hansen (2014) or Chen and Li (2013).

If category-specific predictors are present, like for example in our party choice application, the parameter vector α has to be included in the penalty. For category-specific predictors there is one coefficient for each predictor, so that an ordinary Lasso penalty is appropriate for α and does not introduce issues with respect to the choice of the reference category. For a simple notation, let the overall parameter vector of the model be given by $\theta^T = (\beta_{\bullet 0}^T, \beta_{\bullet 1}^T, \dots, \beta_{\bullet p}^T, \alpha^T)$, where $\beta_{\bullet 0}^T = (\beta_{10}, \dots, \beta_{k-1,0})$ denotes the intercept vector. The parameter θ has length $d = ((p+1)(k-1) + L)$ in a model with an intercept, p global and L category-specific predictors. For this general case, we propose the penalty

$$J(\theta) = \psi \sum_{j=1}^p \phi_j \|\beta_{\bullet j}\| + (1 - \psi) \sum_{l=1}^L \varphi_l |\alpha_l|, \quad (5)$$

where ψ is an additional tuning parameter that balances the penalty on the global and the category-specific variables. Unless stated otherwise, we always use $\psi = 0.5$. The parameters ϕ_j and φ_l are weights that assign different amounts of penalization to different parameter groups. Following the arguments of Yuan and Lin (2006), $\phi_j = \sqrt{k-1}$ and $\varphi_l = 1$ are used as a default, which guarantees that it is “fair” to use the same λ for both the parameter groups $\beta_{\bullet j}$ and the single parameters α_l despite their different size.

The penalty in (5), which contains penalty (4) as a special case, explicitly accounts for the complex categorical structure of the general multinomial logit model. Therefore, we call the corresponding penalized estimator *Categorically Structured Lasso* (CATS Lasso).

2.3. Categorical predictors

The penalty (5), which enforces variable selection in multinomial logit models, was constructed under the restriction that all predictors are metric—more generally, it is suitable for all predictors that would enter a standard univariate GLM with one degree of freedom. This includes metric predictors and binary ones with dummy coding. For the more general case in which a predictor x_j enters a GLM with p_j parameters, extensions of our penalty are needed. The most prominent example for this situation are categorical predictors with $p_j + 1$ categories which typically enter a GLM through p_j dummy variables.

For this case, Yuan and Lin (2006) argued that a grouped penalization and thus grouped selection is more suitable than the selection of single coefficients performed by the ordinary Lasso. For this purpose, they proposed the Group Lasso, which was extended to GLMs in Meier et al. (2008). If such multi-category predictors are to be included in a multinomial logit model one has a vector of parameters for each response category and each predictor, that is, for response category r and variable j , one has $\beta_{rj\bullet}^T = (\beta_{rj1}, \dots, \beta_{rjp_j})$. Thus, there are two “layers” of grouping: First, all coefficients that belong to one actual predictor should enter or leave the model jointly, that is, the whole vector $\beta_{rj\bullet}$ should be set to zero if x_j is found to be irrelevant for category r . Second, the underlying principle of the categorically structured approach says that all the influences $\beta_{rj\bullet}$, $r = 1, \dots, k-1$ of x_j on the different response categories should be penalized in a grouped way. This means that the vector $\beta_{j\bullet}^T = (\beta_{1j\bullet}^T, \dots, \beta_{k-1,j\bullet}^T)$ of all coefficients that are linked to x_j should be treated as one big parameter group. Hence, for multi-category response models the concept of the Group Lasso for univariate models has to be combined with the Categorically Structured Lasso. This also emphasizes that the traditional Group Lasso and CATS Lasso are conceptually different from one another, despite their mathematical similarity.

The same argument also applies to the vector of the category-specific variables α ; if the effect of variable l is determined by the parameter vector $\alpha_l^T = (\alpha_{l1}, \dots, \alpha_{lp_l})$, the whole vector is penalized.

Including categorical predictors yields the CATS Lasso in its most general form:

$$J(\boldsymbol{\theta}) = \psi \sum_{j=1}^p \phi_j \|\boldsymbol{\beta}_{\bullet j}\| + (1 - \psi) \sum_{l=1}^L \varphi_l \|\boldsymbol{\alpha}_{l\bullet}\|, \quad (6)$$

where the weights ϕ_j and φ_l now default to $\phi_j = \sqrt{(k-1)p_j}$ and $\varphi_l = \sqrt{p_l}$. For notational simplicity, we consider metric or binary predictors for the rest of this paper and focus our exposition on the penalty from (5). However, software for the general case is available and has been used in the application.

2.4. Improved variable selection

Like Lasso, Group Lasso, Fused Lasso and other estimators with sparsity-inducing penalties, CATS Lasso is biased. In particular for large values of λ , the sparsity of solutions comes with biased estimates of selected variables. Since a suitable value of λ is usually unknown, it is typically chosen by optimizing some appropriate criterion as, for example, cross-validation or AIC. In practice, however, one can frequently observe that a tuning parameter is chosen for which some weak predictors are not removed. In general, it is desirable to apply a high degree of penalization to weak predictors and mild penalization to strong ones. Two approaches to achieve this goal are discussed next.

2.4.1. Adaptive CATS Lasso

For the ordinary Lasso, Zou (2006) showed that the induced variable selection is inconsistent in certain scenarios and offered a remedy called adaptive Lasso. The same issue and a similar solution were discussed for the Group Lasso by Wang and Leng (2008). The CATS penalty from (5) will suffer from the same problem if simple weights $\phi_j = \sqrt{k-1}$ and $\varphi_l = 1$ are employed. Therefore, we use the idea of adaptive penalties to obtain adaptive CATS Lasso, in which the weights ϕ_j and φ_l are replaced by

$$\phi_j^a = \frac{\sqrt{k-1}}{\|\hat{\boldsymbol{\beta}}_{\bullet j}^{\text{ML}}\|}, \quad \varphi_l^a = \frac{1}{|\hat{\alpha}_l^{\text{ML}}|}, \quad (7)$$

where $\hat{\boldsymbol{\beta}}_{\bullet j}^{\text{ML}}$ and $\hat{\alpha}_l^{\text{ML}}$ denote the respective ML estimates. The basic idea of this adaptive penalization is that the norm of ML estimates of parameter groups belonging to irrelevant predictors asymptotically converges to zero, yielding a strong penalization, while relevant predictors are penalized less severely. In the simulation studies in Section 4, we demonstrate that these adaptive weights improve the quality of both variable selection and predictive performance of CATS Lasso. It should be noted that the ML estimates can be replaced by arbitrary \sqrt{n} -consistent estimates, for example an asymptotically vanishing ridge penalty. Such a ridge penalty is used in our implementation whenever ML estimates do not exist, for example in the $d > n$ case. More technical details are given in the Appendix.

2.4.2. Refitting

A further concept to improve variable selection is to use the penalized estimator for selection purposes only and to perform an unpenalized refit on the set of active predictors, that is, those with non-zero coefficients. This refitting technique was mentioned, for example, in Efron et al. (2004) under the name “Lars-OLS hybrid” and in Candès and Tao (2007) as “Gauss-Dantzig selector”. Our simulation studies show that the variable selection performance of CATS Lasso can improve drastically by such an ML refit. This can be explained by the decoupling of bias and variable selection that is achieved via refitting: If the final estimator is obtained by an unpenalized refit, λ only steers variable selection, so that it can be chosen as large as necessary without having to worry about bias. If refitting is used, one therefore typically observes that higher values of λ are chosen than for estimators without refit.

3. Estimation

For the computation of the estimator proposed in the previous section, we consider maximization of the general penalized log-likelihood given by $l_p(\boldsymbol{\theta}) = l(\boldsymbol{\theta}) - \lambda J(\boldsymbol{\theta})$ with a d -dimensional parameter vector $\boldsymbol{\theta}$, a concave and continuously differentiable log-likelihood $l(\boldsymbol{\theta})$ and a convex penalty term $J(\boldsymbol{\theta})$. The penalized maximum likelihood estimator is defined by

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\operatorname{argmax}} l_p(\boldsymbol{\theta}) = \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\operatorname{argmin}} (-l(\boldsymbol{\theta}) + \lambda J(\boldsymbol{\theta})). \quad (8)$$

CATS Lasso estimates are obtained as a special case if the log-likelihood $l(\boldsymbol{\theta})$ is that of multinomial logit models and the penalty term $J(\boldsymbol{\theta})$ has the form given in (5). In the following, we briefly show how this general maximization problem can be solved with the fast iterative shrinkage-thresholding algorithm (FISTA) of Beck and Teboulle (2009). Technically speaking, FISTA belongs to the class of proximal gradient methods in which only the log-likelihood and its gradient, but no higher-order derivatives are used.

This is particularly useful for multinomial models because they are inherently higher dimensional than univariate GLMs of comparable size. Consider a multinomial logit model with $p = 100$ and $k = 10$. Although this model is not excessively large compared to, for example, gene expression data where $p > 10\,000$ is quite common, the corresponding Fisher matrix would be of size 1000×1000 . Working with or storing such a matrix, as is required for the traditional Fisher scoring method, is very costly and outright impossible for large models. Although the Fisher matrix is avoided and the problem is non-smooth, FISTA achieves quadratic convergence, which is known to be optimal among black-box first-order methods for smooth convex optimization (Nemirovski, 1994; Nesterov, 2004). Thus, it combines quick convergence with cheap iterates that are well-suited for the specific challenges of multinomial logit models. In particular, it turned out to be faster than the block coordinate descent algorithm of Meier et al. (2008) which we tested first, exploiting a complicated representation of multinomial logit models as multivariate GLMs.

Technical details of the algorithm and our implementation are found in the [Appendix](#).

3.1. Proximal gradient methods for penalized log-likelihood problems

Let $\nu > 0$ denote an inverse stepsize parameter. If one sets $\lambda = 0$, the update of traditional gradient methods for the (smooth) minimization of $-l(\theta)$ has the following form:

$$\hat{\theta}^{(t+1)} = \hat{\theta}^{(t)} + \frac{1}{\nu} \nabla l(\hat{\theta}^{(t)}). \quad (9)$$

In proximal gradient methods, the non-smoothness of the actual optimization problem (8) is accounted for by use of the so-called proximal operator.

Given a search point $\mathbf{u} \in \mathbb{R}^d$, the proximal operator associated with penalty $J(\theta)$ is defined by

$$\mathcal{P}_\lambda(\mathbf{u}) = \operatorname{argmin}_{\theta \in \mathbb{R}^d} \left(\frac{1}{2} \|\theta - \mathbf{u}\|^2 + \lambda J(\theta) \right).$$

Using the gradient step from (9) as the search point, the basic iteration of proximal gradient methods is given by

$$\hat{\theta}^{(t+1)} = \mathcal{P}_{\frac{\lambda}{\nu}} \left(\hat{\theta}^{(t)} + \frac{1}{\nu} \nabla l(\hat{\theta}^{(t)}) \right). \quad (10)$$

Thus, the building blocks required to solve (8) with a proximal gradient method are formulas for the log-likelihood and the score function of the general multinomial logit model as well as a formula for the computation of the proximal operator associated with the CATS Lasso penalty given in (5). These building blocks are presented next. Algorithmic details about proximal gradient methods and the acceleration scheme used by FISTA are found in the [Appendix](#).

3.2. Log-likelihood, score function and proximal operator

3.2.1. Log-likelihood of multinomial logit models

For each actual observation $y_i \in \{1, \dots, k\}$, we define, for $r = 1, \dots, k-1$, a set of pseudo-observations y_{ir} , given by $y_{ir} = 1$ if $y_i = r$ and $y_{ir} = 0$ otherwise. With linear predictors η_{ir} of the form

$$\eta_{ir} = \beta_{r0} + \mathbf{x}_i^T \boldsymbol{\beta}_r + (\mathbf{w}_{ir} - \mathbf{w}_{ik})^T \boldsymbol{\alpha}, \quad r = 1, \dots, k-1,$$

the log-likelihood can conveniently be written as

$$l(\theta) = \sum_{i=1}^n \left(\sum_{r=1}^{k-1} y_{ir} \eta_{ir} - \log \left(1 + \sum_{s=1}^{k-1} \exp(\eta_{is}) \right) \right). \quad (11)$$

Note that the logit link used in multinomial logit models is the canonical link for the multinomial distribution. Following arguments of Fahrmeir and Kaufmann (1985), this automatically yields concavity of (11).

3.2.2. Score function

We partition the score function, which is defined as the gradient of the log-likelihood, in the following way:

$$\nabla l(\theta)^T = s(\theta)^T = \left(s(\boldsymbol{\beta}_{\bullet 0})^T \mid s(\boldsymbol{\beta}_{\bullet 1})^T \mid \dots \mid s(\boldsymbol{\beta}_{\bullet p})^T \mid s(\boldsymbol{\alpha})^T \right).$$

To be able to give $s(\theta)$ in a concise form, we use the following notation: For $r = 1, \dots, k-1$, let $\mathbf{y}_r^T = (y_{1r}, \dots, y_{nr})$ and $\boldsymbol{\pi}_r^T = (\pi_{1r}, \dots, \pi_{nr})$ denote vectors that pool the observations and estimated probabilities for category r across all

observations. Additionally, let $\mathbf{x}_j^T = (x_{1j}, \dots, x_{nj})$ denote the vector of all observations of the j th predictor. Furthermore, define

$$\mathbf{V}_r = \begin{pmatrix} v_{11r} & \cdots & v_{1Lr} \\ \vdots & \ddots & \vdots \\ v_{n1r} & \cdots & v_{nLr} \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_{k-1} \end{pmatrix} \quad \text{and} \quad \boldsymbol{\pi} = \begin{pmatrix} \boldsymbol{\pi}_1 \\ \vdots \\ \boldsymbol{\pi}_{k-1} \end{pmatrix},$$

where $v_{ilr} = w_{ilr} - w_{ilk}$, $r = 1, \dots, k-1$ denotes the effect with which the category-specific predictors enter the model. The overall model matrix of size $n(k-1) \times ((p+1)(k-1) + L)$ is

$$\mathbf{Z} = \left(\begin{array}{c|c|c|c|c|c} \mathbf{1}_n & & & & & \\ & \ddots & & & & \\ & & \mathbf{1}_n & & & \\ & & & \mathbf{x}_1 & & \\ & & & & \ddots & \\ & & & & & \mathbf{x}_p \\ & & & & & & \ddots & \\ & & & & & & & \mathbf{V}_1 \\ & & & & & & & \vdots \\ & & & & & & & \mathbf{V}_{k-1} \end{array} \right).$$

With these definitions, the score function, partitioned as above, is given by

$$s(\boldsymbol{\theta}) = \mathbf{Z}^T (\mathbf{y} - \boldsymbol{\pi}). \quad (12)$$

3.2.3. Proximal operator

Let $\tilde{\boldsymbol{\theta}} \in \mathbb{R}^d$ denote a generic search point that is equivalently partitioned as the overall parameter vector $\boldsymbol{\theta}$ of the considered multinomial logit model, that is, $\tilde{\boldsymbol{\theta}}^T = (\tilde{\boldsymbol{\beta}}_{\bullet 0}^T, \tilde{\boldsymbol{\beta}}_{\bullet 1}^T, \dots, \tilde{\boldsymbol{\beta}}_{\bullet p}^T, \tilde{\boldsymbol{\alpha}}^T)$. Additionally, define $\lambda_1 = \lambda\psi$ and $\lambda_2 = \lambda(1-\psi)$. Then the CATS-penalty from (5) can be written as

$$\lambda J(\boldsymbol{\theta}) = \lambda_1 \sum_{j=1}^p \phi_j \|\boldsymbol{\beta}_{\bullet j}\| + \lambda_2 \sum_{l=1}^L \varphi_l |\alpha_l|.$$

The projection of the search point $\tilde{\boldsymbol{\theta}}$ on the constraint region belonging to this penalty is then given by the proximal operator

$$\mathcal{P}_\lambda(\tilde{\boldsymbol{\theta}}) = \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\operatorname{argmin}} \left(\frac{1}{2} \|\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}\|^2 + \lambda J(\boldsymbol{\theta}) \right).$$

Due to the block-separability of this specific penalty, the proximal operator can be written as the sum

$$\mathcal{P}_\lambda(\tilde{\boldsymbol{\theta}}) = \sum_{j=0}^p \mathcal{P}_{\lambda_1}(\tilde{\boldsymbol{\beta}}_{\bullet j}) + \sum_{l=1}^L \mathcal{P}_{\lambda_2}(\tilde{\alpha}_l),$$

where

$$\begin{aligned} \mathcal{P}_{\lambda_1}(\tilde{\boldsymbol{\beta}}_{\bullet 0}) &= \underset{\boldsymbol{\beta}_{\bullet 0} \in \mathbb{R}^{k-1}}{\operatorname{argmin}} \frac{1}{2} \|\boldsymbol{\beta}_{\bullet 0} - \tilde{\boldsymbol{\beta}}_{\bullet 0}\|^2 = \tilde{\boldsymbol{\beta}}_{\bullet 0}, \\ \mathcal{P}_{\lambda_1}(\tilde{\boldsymbol{\beta}}_{\bullet j}) &= \underset{\boldsymbol{\beta}_{\bullet j} \in \mathbb{R}^{k-1}}{\operatorname{argmin}} \left(\frac{1}{2} \|\boldsymbol{\beta}_{\bullet j} - \tilde{\boldsymbol{\beta}}_{\bullet j}\|^2 + \lambda_1 \phi_j \|\boldsymbol{\beta}_{\bullet j}\| \right), \quad j = 1, \dots, p, \end{aligned} \quad (13)$$

and

$$\mathcal{P}_{\lambda_2}(\tilde{\alpha}_l) = \underset{\alpha_l \in \mathbb{R}}{\operatorname{argmin}} \left(\frac{1}{2} (\alpha_l - \tilde{\alpha}_l)^2 + \lambda_2 \varphi_l |\alpha_l| \right), \quad l = 1, \dots, L. \quad (14)$$

With $(u)_+ = \max(u, 0)$, the following analytic solutions to (13) and (14) can easily be derived from the Karush–Kuhn–Tucker conditions:

$$\mathcal{P}_{\lambda_1}(\tilde{\boldsymbol{\beta}}_{\bullet j}) = \left(1 - \frac{\lambda_1 \phi_j}{\|\tilde{\boldsymbol{\beta}}_{\bullet j}\|} \right)_+ \tilde{\boldsymbol{\beta}}_{\bullet j}, \quad (15)$$

$$\mathcal{P}_{\lambda_2}(\tilde{\alpha}_l) = \left(1 - \frac{\lambda_2 \varphi_l}{|\tilde{\alpha}_l|} \right)_+ \tilde{\alpha}_l. \quad (16)$$

3.3. Choice of tuning parameters

We choose the main tuning parameter λ based on cross-validation. But in the general case, the additional tuning parameter ψ , which balances the penalization between parameters belonging to global and to category-specific predictors, has to be chosen for given λ . We performed various simulation studies (not shown) in which we simultaneously cross-validated over both λ and ψ . The ψ estimated by cross-validation always remained within the interval $[0.35, 0.65]$. In addition, the difference between the model with cross-validated ψ and the model using the default value of $\psi = 0.5$ was

very minor. It seems that the default values for the weights ϕ_j and φ_l are already balancing out the penalty well enough. For efficiency reasons, we therefore recommend to choose only λ via cross-validation and to keep ψ fixed at 0.5.

For all simulations and the real data example in this paper, we used 10-fold cross-validation based on the deviance. Alternatively, model selection criteria like AIC can be employed, using the effective degrees of freedom given in Yuan and Lin (2006). For CATS Lasso, the formula is

$$\hat{df} = (k - 1) + \sum_{j=1}^p \left(\mathbb{1}(\|\hat{\beta}_{\bullet j}\| > 0) + \frac{\|\hat{\beta}_{\bullet j}\|}{\|\hat{\beta}_{\bullet j}^{\text{ML}}\|} (k - 2) \right) + \sum_{l=1}^L \mathbb{1}(|\hat{\alpha}_l| > 0), \quad (17)$$

where the first term corresponds to the unpenalized intercept vector and $\mathbb{1}(\cdot)$ denotes the indicator function.

4. Simulation study

Before analyzing the German Longitudinal Election Study, we first give the results of a small simulation study which illustrates the performance of our approach. In particular, four different versions of CATS Lasso can be derived by using or not using adaptive weights and/or ML refitting. The simulations demonstrate the performance of these four variants under various settings and provide a guideline on when to use which of the alternative versions. The second aim of the simulation study is to compare our categorically structured approach with the unstructured Lasso. In order to simplify the presentation, we always refer to CATS Lasso as “CATS” and to the ordinary, unstructured Lasso as “Lasso” for the remainder of this section. All simulations were performed in R (R Development Core Team, 2014) using a self-written implementation that can be obtained from the authors.

4.1. Simulation settings

4.1.1. Scenarios

We consider a small and a large model. The small model uses 5 response categories with 4 relevant and 4 noise variables, giving a total of 8 global predictors. Additionally, 4 category-specific predictors are available of which 2 are relevant. The size of this model is roughly equivalent to that of our application. The large model consists of 10 response categories, 60 global and 20 category-specific predictors of which 20 and 8 are relevant, respectively. For each model, the coefficients of the relevant predictors are independently drawn at random from the set $\{-1, -0.5, 0.5, 1, 1.5, 2, 2.5, 3\}$, yielding a true coefficient vector θ^* . The coefficients of noise variables are always zero.

Both global and category-specific predictors were drawn from a multivariate Gaussian with an equi-correlation of 0.2 (small model) or 0.6 (large model) between all predictors. Using these predictors and the true coefficient vector θ^* , the true probabilities for each observation were computed and then used to draw the response from a multinomial distribution.

The small model is tested for $n = 40$ and $n = 200$, the large one for $n = 500$. Including intercepts, the size of the overall parameter vector to be estimated is $d = 40$ for the small model and $d = 569$ for the large model. In the large model, binary predictors were included, but this had only a negligible effect on the simulation results. To test the methods for their behavior under overdispersion, the large simulation was repeated with a response drawn from a Dirichlet-multinomial distribution, using R package `dirmult` (Tvedebrink, 2010) with parameters chosen to obtain 40% overdispersion. For all settings, the true coefficient vector is drawn once, followed by 100 replications of data generation and model estimation. All results reported below are averages (of the respective quantity) across the 100 simulation runs.

4.1.2. Comparison of methods

To compare the methods, their estimation and prediction accuracy as well as variable selection performance are evaluated. First, with $\hat{\theta}^{(i)}$ denoting the estimator for the i th replication, we define the squared error for this replication as $(\hat{\theta}^{(i)} - \theta^*)^T (\hat{\theta}^{(i)} - \theta^*) / d$. From these squared errors, we compute the mean squared error (MSE). Second, prediction accuracy is evaluated by drawing a test set of $n_{\text{test}} = 3n$ new observations from the true model and then computing the predictive deviance on this test set. Third and last, we report the false positive and false negative rates (in terms of variable selection). The false positive rate (FPR) for variable selection is the percentage of noise variables whose coefficient vector is incorrectly estimated as non-zero. The false negative rate (FNR) for variable selection is the percentage of relevant predictors whose estimated coefficient vector is falsely set to zero.

4.2. Results

4.2.1. Small model

The simulation results for the small model are summarized in Fig. 1. The methods compared are the ML estimator, CATS in simple, adaptive (“ada”), refitted (“rf”) and adaptive plus refitted form as well as the unstructured Lasso in the same four variants as CATS. The left column of Fig. 1 gives the results for the case with more observations than parameters, the right column contains results for the $n \leq d$ case. The first row shows squared errors and the second one the predictive deviance

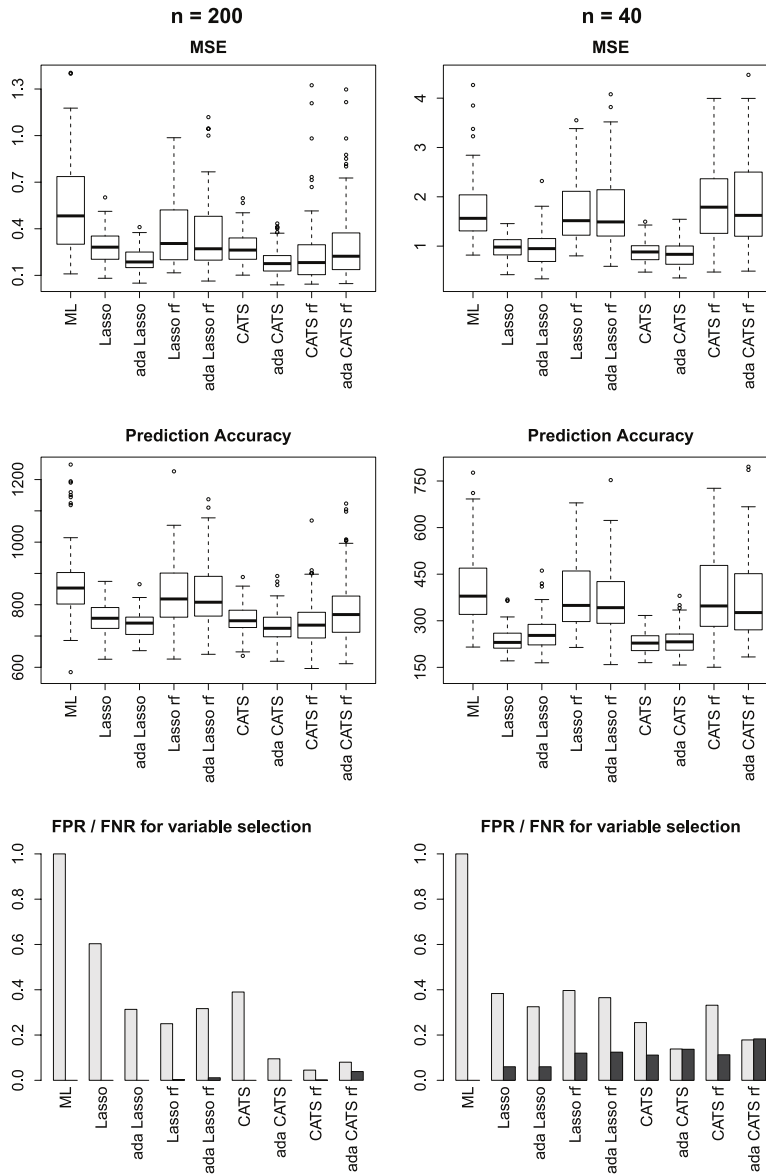


Fig. 1. Simulation results for the small model ($d = 40$): MSE, predictive deviance and False Positive Rate (FPR, gray) and False Negative Rate (FNR, black) for $n = 200$ (left column) and $n = 40$ (right column).

on a test set of, respectively, $n_{\text{test}} = 600$ (left column) and $n_{\text{test}} = 120$ (right column) new observations. Both quantities are given as boxplots to visualize variability and outlier behavior. Extreme values were omitted for better clarity of the plots. In the third row, false positive (gray) and false negative (black) rates are shown. For $n = 200$ (left column), one can clearly see that all regularized estimators perform better than the ML estimator in terms of MSE. Moreover, the large box indicates that the ML estimator is rather instable. For both CATS and Lasso, the adaptive version without refitting performs best in terms of MSE and prediction accuracy. Comparing each of the four variants between CATS and Lasso, CATS always comes out slightly ahead. However, the biggest advantage of CATS over Lasso is visible in the false positive and false negative rates, shown in the bottom left plot.

While plain CATS does not select satisfactorily, both refitting and adaptive weights improve the variable selection of CATS by a large margin, confirming the issues of simple CATS that were discussed in Section 2.4. Lasso shows substantially worse variable selection behavior that, in contrast to CATS, is not mitigated by refitting or adaptive weights.

For $n = 40$ (right column), a model is fitted with as many parameters as observations. In this case, “true” ML estimates do not exist anymore, so we added a very small ridge penalty whenever an unpenalized fit for this model was required. This is the case for the ML estimator as well as all methods with refit. It can immediately be seen from Fig. 1 that for $n = 40$, all these methods in which the final model is obtained by an unpenalized ML fit perform significantly worse.

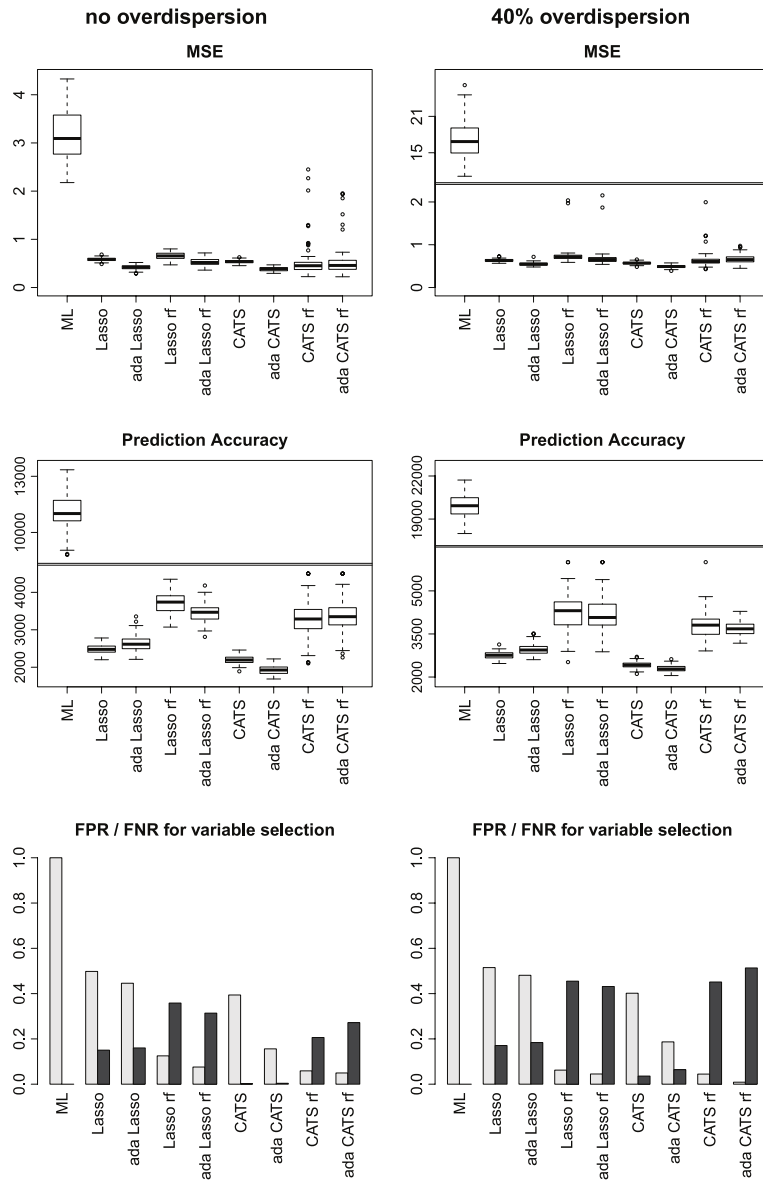


Fig. 2. Simulation results for the large model ($d = 569$): MSE, predictive deviance and False Positive Rate (FPR, gray) and False Negative Rate (FNR, black) for $n = 500$ and no (left column) and 40% (right column) overdispersion.

In terms of MSE and prediction accuracy, differences between Lasso, CATS and their respective adaptive versions are minor. When it comes to variable selection, however, adaptive CATS shows by far the best performance of all considered methods and includes only half as many noise variables as the best version of Lasso. It is noteworthy that the false negative rates are much higher throughout all methods than for $n = 200$.

4.2.2. Large model

The results of the simulations for the large model are shown in Fig. 2. The left column of Fig. 2 shows the results for the large model with $n = 500$, that is, for a setting in which significantly less observations than parameters to be estimated are available. Once again, a small ridge penalty is added to the ML estimator and during refitting since an unregularized estimator does not exist in this scenario.

In terms of MSE, the differences between the regularized estimators are not very pronounced while the ML estimator – as expected – shows an extremely poor performance. To show the predictive deviance on $n_{\text{test}} = 1500$ observations the axis had to be broken to include the ML estimator. The plot also shows that refitting techniques, in which the final model is obtained by an unpenalized ML refit on the active set of the regularized estimator, perform substantially worse. But it is essential that adaptive CATS distinctly outperforms Lasso. The plot on the bottom left of Fig. 2 shows that all

methods with refitting have high false negative rates. Moreover, all versions of Lasso have both higher false positive and false negative rates than their corresponding counterpart of CATS. For the considered scenario, adaptive CATS offers the best variable selection properties by a large margin. The right column of Fig. 2 shows the results for the large model with 40% overdispersion. Compared to the same setting without overdispersion, the performance of the ML estimator (equipped with a small ridge penalty) deteriorates even more. A comparison of the left and right columns in Fig. 2 shows that the impact of overdispersion on the MSE of penalized estimators is minor, but the predictive performance of all methods suffers in the presence of overdispersion.

For example, the mean predictive deviance of adaptive CATS increases from 1759 to 2275, or about 29.3%. The bottom right plot shows that false positive rates are mostly unaffected by overdispersion, but the false negative rates increase, especially for estimators that include refitting. Once again, adaptive CATS shows the best variable selection properties and seems to be rather stable under overdispersion.

4.2.3. Summary

To sum up the simulation study, regularized estimators offer more accurate estimates and prediction than the unpenalized ML estimator—even if substantially more observations than parameters to be estimated are available. CATS outperforms the ordinary Lasso in all settings and in terms of all criteria that were used to compare them. The degree by which CATS outperforms Lasso grows with the number of response categories, in particular in terms of variable selection. Comparing the simple, adaptive and refitted version of CATS, adaptive CATS is clearly recommended if few observations are available, relative to the size of the model. Our simulation results indicate that the adaptive version of CATS (and also Lasso) is less affected by an instability of this estimator than refitting approaches.

In data-rich situations, refitting performed marginally better than adaptive weights when it comes to variable selection, but slightly worse in terms of MSE and prediction accuracy. The combination of adaptive weights and refitting did not offer an advantage, but was more unstable in several cases and is therefore not recommended.

From the boxplots in the left column of Fig. 1, one can see that the refitting approach does not only perform worse than the adaptive one in terms of MSE and prediction but also has a larger variability and produces more outliers. Owing to higher stability, we therefore prefer adaptive CATS over refitted CATS even in $n > d$ settings.

5. Regularized analysis of party choice in Germany

5.1. Data description

The data we consider come from the German Longitudinal Election Study. We focus on a dataset that contains the party on which 816 study participants intended to vote during the 2009 election for the German parliament, the Bundestag. To be able to explain their party choice behavior, nine individual characteristics of the voters are available. These global predictors are *gender* (1: male, 0: female), *regional provenance* (*west*; 1: former West Germany, 0: former East Germany), *age* (mean-centered), *union* (1: membership in a union, 0: otherwise), *high school degree* (1: yes, 0: no), *unemployment* (1: currently unemployed, 0: otherwise), *political interest* (1: less interested, 0: very interested), *satisfaction with the functioning of democracy* (*democracy*; 1: not satisfied, 0: satisfied) and *religion* (0: Protestant, 1: Catholic, 2: otherwise).

Additionally, the interviewed persons rated the position of the parties on political issues on a scale from 1 to 11. They also rated their own position on these topics on the same scale. From this information, the distance between the voters' own position (the "ideal point") and the perceived position of the party was computed, that is, the absolute value of the difference of the two variables. The resulting distances take values between 0 and 10 and can be considered measures of agreement between voter and party and are used as category-specific predictors. This approach has strong connections to spatial election theory which assumes that each party is characterized by a position in a finite-dimensional space, with each dimension corresponding to one political issue. Spatial election theory then explains the party choice of voters by a utility function that depends on the distance between the voter's own position and that of the parties within the space of policy-dimensions. For further details on spatial election theory, see, for example, Thurner and Eymann (2000). Here, four political issues were considered: *taxes* ("do you prefer low taxes and few public spending or high taxes with lots of public spending?"), the attitude toward *immigration* and *nuclear energy* as well as the positioning on a political left--right scale.

Correlations between voter-specific covariates and the party-specific issue variables are very small. Among global predictors, only the correlations between *west* and *religion* (Catholic: 0.34, Other: -0.36) are of magnitude larger than 0.25. The correlation among issue variables is higher and ranges from 0.33 to 0.44. Overall, the correlation in the data is negligible or weak for all global predictors and moderate among the category-specific predictors.

5.2. Results and interpretation

To model the party choice behavior of voters in Germany, a multinomial logit model with the chosen party as response variable is used with the explanatory variables described in the previous paragraph. Of the five available parties, the CDU was chosen as reference. We fitted this model using adaptive CATS Lasso, that is, a penalized multinomial logit model with

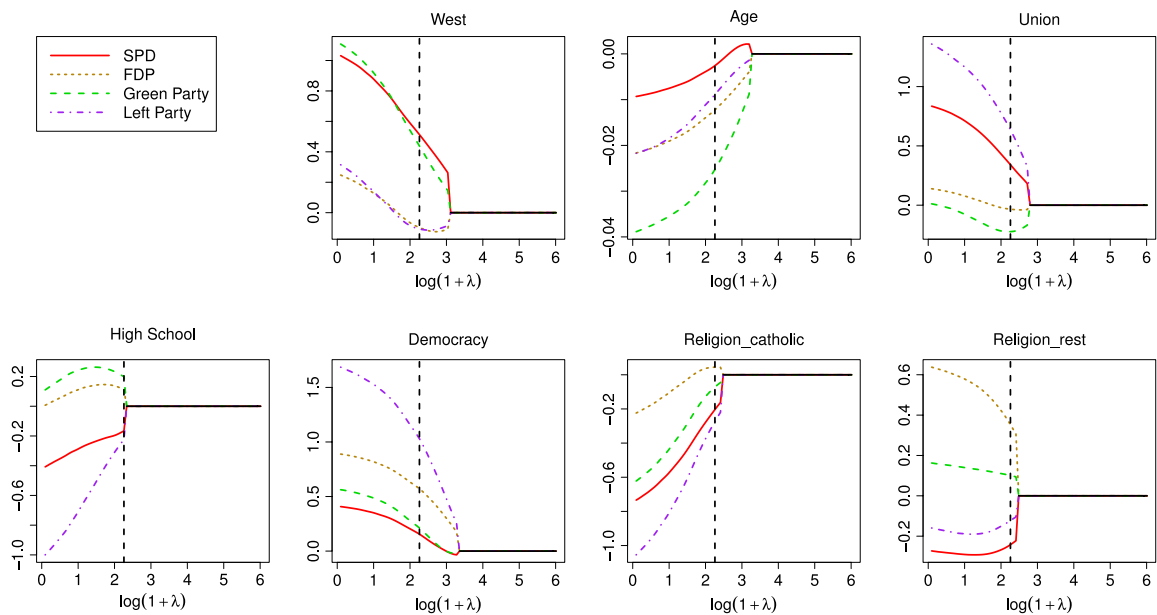


Fig. 3. Coefficient builds for the selected global variables of party choice data.

Table 1
Estimated regression coefficients for the global predictors of the party choice data. Numbers in parentheses show bootstrapped selection probabilities.

| | SPD | | FDP | | Green party | | Left party | |
|------------------|--------|---------|--------|---------|-------------|---------|------------|---------|
| Intercept | −0.890 | (1.000) | −1.786 | (1.000) | −1.753 | (1.000) | −1.340 | (1.000) |
| Gender | 0 | (0.291) | 0 | (0.291) | 0 | (0.291) | 0 | (0.291) |
| West | 0.738 | (0.998) | 0.025 | (0.998) | 0.735 | (0.998) | 0.006 | (0.998) |
| Age | −0.007 | (1.000) | −0.017 | (1.000) | −0.034 | (1.000) | −0.015 | (1.000) |
| Union | 0.499 | (0.997) | 0.006 | (0.997) | −0.202 | (0.997) | 0.940 | (0.997) |
| High School | −0.212 | (0.839) | 0.144 | (0.839) | 0.229 | (0.839) | −0.320 | (0.839) |
| Unemployment | 0 | (0.000) | 0 | (0.000) | 0 | (0.000) | 0 | (0.000) |
| Pol. Interest | 0 | (0.028) | 0 | (0.028) | 0 | (0.028) | 0 | (0.028) |
| Democracy | 0.332 | (1.000) | 0.756 | (1.000) | 0.431 | (1.000) | 1.468 | (1.000) |
| ReligionCatholic | −0.427 | (0.998) | −0.027 | (0.998) | −0.282 | (0.998) | −0.601 | (0.998) |
| ReligionRest | −0.303 | (0.998) | 0.520 | (0.998) | 0.129 | (0.998) | −0.205 | (0.998) |

the penalty from (5) and the weights given in (7). The only exception is the religion of the voter which enters the model with two dummies, so that the advanced methodology from Section 2.3 is used to jointly penalize both dummies.

Fig. 3 shows the resulting coefficient paths for those variables that were selected by CATS Lasso. They show how the different coefficients change if the tuning parameter λ is varied with λ plotted on a logarithmic scale. The left end of the coefficient paths corresponds to zero penalization and thus shows the ML estimator. The vertical line marks the value of λ that was chosen by 10-fold cross-validation. The global predictors that were removed from the model are gender, unemployment and political interest. One can immediately see the structured selection performed by CATS Lasso. While the variable high school was barely selected, removing democracy, age, west and union would require substantially more penalization. The estimated coefficients for the global variables as well as selection probabilities over $B = 1000$ bootstrap samples are given in Table 1.

Relative to the other three parties, the SPD and the Green Party are performing much better in former West Germany than in former East Germany, given all other variables remain fixed. The coefficients of age are negative for all four parties, which means that the reference party, the CDU, is strongest among elderly people. Growing age is most detrimental to the Green Party, which matches expectations as this party was founded in 1980, so that many older people were already adults with a developed political preference before this party even existed. Voters that are members in a labor union are distinctly more likely to vote for the SPD or the Left Party, which again is expected as these parties are traditionally strong among blue-collar workers. Voters with a high-school degree are more likely to vote in favor of the FDP or the Green Party, which focus on liberal and environmental topics, respectively, whereas the SPD and the Left Party have more success among voters without a high-school degree. Not being satisfied with the functioning of democracy decreases the chances of voting for the CDU, which is not surprising as the CDU is the most conservative party. The Left Party, on the other hand, benefits strongly if the voter is discontent with democracy. Therefore, the Left Party and, to a lesser degree, the FDP can be considered the parties of choice for “protest voters”. Catholic people prefer the FDP or the CDU and are noticeably less likely to vote for the

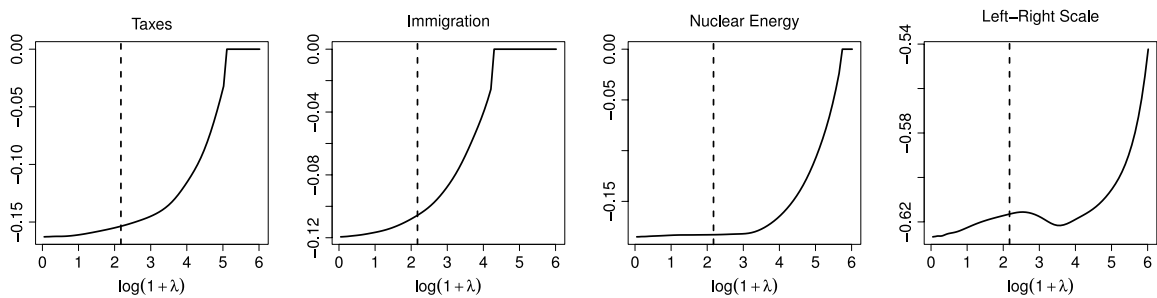


Fig. 4. Coefficient buildups for category-specific variables of party choice data.

Table 2

Performance criteria for ML, CATS and Lasso models for different reference categories.

| Estimator | Reference | Deviance | AIC | BIC | Nonzero coefs | Pred. Dev. |
|-----------|-----------|----------|---------|---------|---------------|------------|
| ML | Any | 1596.27 | 1692.27 | 1918.08 | 48 | 457.62 |
| CATS | CDU | 1614.13 | 1671.64 | 1806.91 | 36 | 452.07 |
| | SPD | 1612.16 | 1670.40 | 1807.37 | 36 | 452.17 |
| | FDP | 1616.66 | 1673.62 | 1807.60 | 36 | 453.19 |
| | Green P. | 1618.46 | 1675.37 | 1809.21 | 36 | 451.46 |
| | Left P. | 1622.84 | 1674.12 | 1794.75 | 36 | 454.10 |
| Lasso | CDU | 1602.08 | 1671.46 | 1834.66 | 36 | 452.83 |
| | SPD | 1599.36 | 1676.74 | 1858.73 | 40 | 456.57 |
| | FDP | 1608.44 | 1667.63 | 1806.84 | 31 | 454.22 |
| | Green P. | 1599.74 | 1671.21 | 1839.30 | 37 | 453.71 |
| | Left P. | 1607.01 | 1677.95 | 1844.64 | 37 | 455.52 |

SPD, Green Party or Left Party. The success of the CDU with Catholic voters was expected as it is the most conservative among the major parties in Germany and defending Christian values is a core part of the party agenda. Being neither Protestant nor Catholic increases the likelihood of voting for the FDP or the Green Party.

To interpret the effect of the category-specific predictors, it is helpful to recall the structure of the linear predictor for party r , which has the form $\eta_{ir} = \beta_{r0} + \mathbf{x}_i^T \boldsymbol{\beta}_r + (\mathbf{w}_{ir} - \mathbf{w}_{ik})^T \boldsymbol{\alpha}$, where the vector \mathbf{x}_i collects the nine global predictors that characterize the i th voter. The vector \mathbf{w}_{ir} denotes the four-dimensional vector of distances between the positioning of the voter and the perceived position of party r for the issues taxes, immigration, nuclear energy and the left–right scale. As a reference point, \mathbf{w}_{ik} contains the distances on these issues between the voter and the CDU. Hence, the variable w_{isl} , $s = 1, \dots, k$ measures the disagreement between voter i and party s on the l th considered political topic, that is, low values correspond to high agreement and vice versa. The coefficient α_l is the effect of the quantity $w_{irl} - w_{ikl}$, which corresponds to the difference between the disagreement with the r th party and the reference party. If this difference is positive, this means that the voter shows stronger agreement with the CDU than with party r on the topic at hand, so that the odds of voting for this party instead of the CDU should decrease. Therefore, we expect the coefficients for the political issues to be negative.

All four political issues were selected by our method. The estimated parameters for the four issues taxes, immigration, nuclear energy and the left–right scale are given by $\hat{\boldsymbol{\alpha}}^T = (-0.160, -0.113, -0.185, -0.622)$, respectively. These estimators all have a selection probability of 1. The signs are negative as expected, meaning that the odds of voting for a party increase if the agreement between a voter and this party on a particular issue is higher than the agreement with the reference party. The corresponding coefficient paths are shown in Fig. 4. The vertical line shows the λ chosen by cross-validation. One can see that the left–right scale has the largest influence on party choice, but it is also the most general of the four considered issues. Comparing Figs. 3 and 4 and keeping in mind that a logarithmic scale is used on the x -axis, one sees that the voter-specific predictors are shrunk to zero much earlier, so that the agreement on political issues can be considered a stronger predictor for party preference.

5.3. Comparison of model performance

To assess the performance of CATS, the deviance, AIC, BIC as well as the number of nonzero coefficients are considered. These quantities are shown in Table 2 for the ML estimator, CATS, the ordinary Lasso and all five possible choices of the reference category. In addition, the out-of-sample performance is examined by random splitting, in which the dataset is randomly split into 600 training observations and 216 test observations. The models are estimated on the training data (including tuning parameter selection) and these estimated models are then used to compute the predictive deviance on the test data. To eliminate the randomness in data splitting, the last column of Table 2 reports the average predictive deviances over 1000 replications of this procedure. Their standard errors range from 0.66 to 0.70.

As expected, the ML estimator provides the closest fit to the data, but is distinctly outperformed by the regularized estimators in terms of AIC, BIC, sparsity and prediction. The Lasso tends to provide models with a slightly closer fit to the data than CATS if the number of parameters is not taken into account, but CATS shows better or comparable BIC values than Lasso for all reference categories while a comparison by AIC is rather inconclusive. Moreover, CATS provides better out-of-sample prediction. It should be noted that the effective degrees of freedom for CATS are computed via formula (17) which depends not only on the number of selected variables but also on the amount of shrinkage. Therefore, the differences between AIC (and BIC) values for CATS for different reference categories do not directly reflect the differences in deviance. By contrast to CATS, the sparsity of Lasso solutions, which is computed following Zou et al. (2007) and Yuan and Lin (2006), indeed depends on the reference category. The weak dependence of CATS on the choice of the reference category was already mentioned in Section 2. It is seen that CATS yields the same number of coefficients for all choices whereas the lasso shows distinct variation. Moreover, the differences between the maximal and minimal values in terms of AIC and BIC are much smaller for CATS than for the Lasso.

6. Concluding remarks

In this paper, we analyze the party choice preference of German voters using a regularized multinomial logit model. Since multinomial models are designed for multi-category responses, they are inherently multivariate. To perform effective variable selection, we adopt a contemporary grouped penalization approach with groups formed by all coefficients that belong to the same predictor variable. In addition to ordinary global predictors like age or gender of the voter, the party choice application contains covariates that vary across parties, like, e.g., stance on immigration and also predictors like religion which are represented by several dummy variables. To regularize such a model, we propose Categorically Structured Lasso (CATS Lasso), a penalization scheme that accounts for the complex model structure arising from this combination of a multivariate response with these different predictor types.

To improve the variable selection and predictive performance, we suggest an adaptive version of CATS Lasso as well as a refitting procedure. Our R implementation uses an efficient algorithm and is the first to support the combination of regularized estimation and category-specific predictors in multinomial logit models.

In a simulation study, we have shown that CATS Lasso outperforms alternative regularization approaches for multinomial models in small and large as well as sparse and data-rich models. The application on the modeling of party choice in Germany demonstrates the success of CATS Lasso on real-world problems. It shows that in the 2009 German parliamentary election, the employment status, gender and political interest of a voter did not influence his or her party preference given all other predictors.

Acknowledgments

This work was partially supported by DFG (German Science Foundation) grant TU 62/5-1 (“Regularisierung für diskrete Datenstrukturen”). We sincerely thank the associate editor and three anonymous reviewers for their helpful comments that allowed us to improve the paper.

Appendix. Algorithmic details

A.1. The fast iterative shrinkage thresholding algorithm (FISTA)

With $\nu > 0$ denoting an inverse stepsize parameter, the iterations of proximal gradient methods for maximizing a penalized log-likelihood $l_p(\theta) = l(\theta) - \lambda J(\theta)$ are defined by Beck and Teboulle (2009)

$$\hat{\theta}^{(t+1)} = \operatorname{argmin}_{\theta \in \mathbb{R}^d} \left(Q_\nu(\theta, \hat{\theta}^{(t)}) := -l(\hat{\theta}^{(t)}) - \nabla l(\hat{\theta}^{(t)})^T (\theta - \hat{\theta}^{(t)}) + \frac{\nu}{2} \|\theta - \hat{\theta}^{(t)}\|^2 + \lambda J(\theta) \right), \quad (\text{A.1})$$

which consists of a linear approximation of the negative log-likelihood at the current value $\hat{\theta}^{(t)}$, a proximity term and the penalty. If one sets $\lambda = 0$, the scheme in (A.1) yields a sequence of unpenalized estimators, denoted here by $\tilde{\theta}^{(t)}$, for which the following explicit form is available:

$$\tilde{\theta}^{(t+1)} = \tilde{\theta}^{(t)} + \frac{1}{\nu} \nabla l(\tilde{\theta}^{(t)}). \quad (\text{A.2})$$

This is the standard formula for the iterates of gradient methods for smooth optimization and intuitively appealing because of the interpretation of the gradient as the direction of steepest ascent. Note that the sequence $\{\tilde{\theta}^{(t)}\}$ converges to the ML estimator. Hence, the update in (A.2) can be considered a one-step approximation to the ML estimator based on the current iterate.

However, we require solutions to (A.1) with an active penalty. To obtain these solutions it is helpful to consider a more explicit and tractable formulation of (A.1) that can be derived using standard convex optimization theory: Via Lagrange duality (Bertsekas et al., 2003), Eq. (8) can equivalently be expressed by

$$\hat{\boldsymbol{\theta}} = \operatorname{argmin}_{\boldsymbol{\theta} \in \mathcal{C}} (-l(\boldsymbol{\theta})), \quad (\text{A.3})$$

where $\mathcal{C} = \{\boldsymbol{\theta} \in \mathbb{R}^d \mid J(\boldsymbol{\theta}) \leq \kappa(\lambda)\}$ is the constraint region corresponding to $J(\boldsymbol{\theta})$ and $\kappa(\lambda)$ is a tuning parameter that is linked to λ by a one-to-one mapping. Given a search point $\mathbf{u} \in \mathbb{R}^d$, the so-called proximal operator associated with $J(\boldsymbol{\theta})$ is defined by

$$\mathcal{P}_\lambda(\mathbf{u}) = \operatorname{argmin}_{\boldsymbol{\theta} \in \mathbb{R}^d} \left(\frac{1}{2} \|\boldsymbol{\theta} - \mathbf{u}\|^2 + \lambda J(\boldsymbol{\theta}) \right) \quad (\text{A.4})$$

and yields the projection of \mathbf{u} onto \mathcal{C} , that is, for any $\mathbf{u} \in \mathbb{R}^d$, one has $\mathcal{P}_\lambda(\mathbf{u}) \in \mathcal{C}$. Using this notation and simple algebra, the proximal gradient iterates defined in (A.1) can also be expressed by the projection

$$\hat{\boldsymbol{\theta}}^{(t+1)} = \mathcal{P}_{\frac{\lambda}{\nu}} \left(\hat{\boldsymbol{\theta}}^{(t)} + \frac{1}{\nu} \nabla l(\hat{\boldsymbol{\theta}}^{(t)}) \right). \quad (\text{A.5})$$

This means that a gradient step toward the ML estimator is performed, starting from the current iterate $\hat{\boldsymbol{\theta}}^{(t)}$, to create a search point, which is then projected on the penalty region by the proximal operator. The analysis in Beck and Teboulle (2009), however, showed that this procedure does not reach the optimal convergence rate and can be accelerated. To achieve this speed-up, the direction given by the line between the estimators in the current and the previous iteration is extrapolated with the help of deliberately chosen acceleration factors a_t :

$$\hat{\boldsymbol{\vartheta}}^{(t)} = \hat{\boldsymbol{\theta}}^{(t)} + \frac{a_{t-1} - 1}{a_t} (\hat{\boldsymbol{\theta}}^{(t)} - \hat{\boldsymbol{\theta}}^{(t-1)}).$$

Instead of the current iterate $\hat{\boldsymbol{\theta}}^{(t)}$, this extrapolated point $\hat{\boldsymbol{\vartheta}}^{(t)}$ is used as the starting point from which the step toward the ML estimator is taken, yielding an “extrapolated search point” which is then projected on the penalty region. Because in practice one does not know a suitable value for the inverse stepsize parameter ν , it is determined by a standard backtracking line search (see, e.g., Bertsekas et al., 2003). In detail, the FISTA method for minimizing (8) is given by

FISTA for computing maximum penalized log-likelihood estimates

(0) **Input:** Formulas for $l(\boldsymbol{\theta})$, $\nabla l(\boldsymbol{\theta})$, $J(\boldsymbol{\theta})$, $\mathcal{P}_\lambda(\cdot)$; starting values $\hat{\boldsymbol{\theta}}^{(0)}$, $\nu_0 > 0$; a tolerance ε .

(1) **Initialize:** Set $\hat{\boldsymbol{\theta}}^{(1)} = \hat{\boldsymbol{\theta}}^{(0)}$, $a_0 = 0$, $a_1 = 1$ and $t = 1$.

(2) **Extrapolate:** $\hat{\boldsymbol{\vartheta}}^{(t)} = \hat{\boldsymbol{\theta}}^{(t)} + \frac{a_{t-1} - 1}{a_t} (\hat{\boldsymbol{\theta}}^{(t)} - \hat{\boldsymbol{\theta}}^{(t-1)})$.

(3) **Line search:**

(i) **Initialize:** $\nu = \nu_{t-1}$.

(ii) **Projection:** $\hat{\boldsymbol{\theta}}^{(t+1)} = \mathcal{P}_{\frac{\lambda}{\nu}} \left(\hat{\boldsymbol{\vartheta}}^{(t)} + \frac{1}{\nu} \nabla l(\hat{\boldsymbol{\vartheta}}^{(t)}) \right)$.

(iii) **Armijo rule:** End if

$-l_p(\hat{\boldsymbol{\theta}}^{(t+1)}) \leq Q_\nu(\hat{\boldsymbol{\theta}}^{(t+1)}, \hat{\boldsymbol{\vartheta}}^{(t)}),$
else multiply ν by 2 and repeat.

(4) **Update:**

$$a_{t+1} = \frac{1 + \sqrt{1 + 4a_t^2}}{2},$$

$$\nu_t = \nu,$$

$$\hat{\boldsymbol{\theta}}^{(t+1)} = \mathcal{P}_{\frac{\lambda}{\nu_t}} \left(\hat{\boldsymbol{\vartheta}}^{(t)} + \frac{1}{\nu_t} \nabla l(\hat{\boldsymbol{\vartheta}}^{(t)}) \right).$$

(5) **Convergence check:** Stop if

$$\frac{\|l_p(\hat{\boldsymbol{\theta}}^{(t+1)}) - l_p(\hat{\boldsymbol{\theta}}^{(t)})\|}{\|l_p(\hat{\boldsymbol{\theta}}^{(t)})\|} \leq \varepsilon,$$

else increase t by 1 and repeat steps (2)–(5).

A.2. Stabilizing ML estimates

ML estimates of the full model are required to compare our regularized estimator with the unpenalized model and to compute the adaptive weights from (7). However, we frequently observed a divergence of parameters to $\pm\infty$ even in the

$d < n$ case. Roughly speaking, this happens when the data space allows for a perfect separation of the response categories. This is always the case for $d > n$, but can also happen if the number of observations per parameter to be estimated is above 1, but rather small. The same problems were also observed and discussed by Friedman et al. (2010).

It is obvious that any kind of p -norm penalization of the parameters (with $p > 0$) will prevent this problem. Therefore, we propose to add a ridge penalty (Hoerl and Kennard, 1970) with a very small tuning parameter (say between 0.01 and 0.1) if coefficients are detected to be diverging:

$$J_R(\boldsymbol{\theta}) = \frac{\lambda_R}{2} \|\boldsymbol{\theta}\|_2^2.$$

The proximal operator associated with this ridge penalty is given by

$$\mathcal{P}_{\lambda_R}(\tilde{\boldsymbol{\theta}}) = \frac{\tilde{\boldsymbol{\theta}}}{1 + \lambda_R}.$$

References

- Argyriou, A., Evgeniou, T., Pontil, M., 2007. Multi-task feature learning. In: Schölkopf, B., Platt, J., Hoffman, T. (Eds.), *Advances in Neural Information Processing Systems*, Vol. 19. MIT Press, pp. 41–48.
- Beck, A., Teboulle, M., 2009. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sci.* 2 (1), 183–202.
- Bertsekas, D., Nedić, A., Ozdaglar, A., 2003. *Convex Analysis and Optimization*. Athena Scientific, Belmont.
- Bühlmann, P., Hothorn, T., 2007. Boosting algorithms: regularization, prediction and model fitting. *Statist. Sci.* 22 (4), 477–505.
- Bühlmann, P., Yu, B., 2003. Boosting with the L2 loss: regression and classification. *J. Amer. Statist. Assoc.* 98 (462), 324–339.
- Candes, E., Tao, T., 2007. The dantzig selector: statistical estimation when p is much larger than n . *Ann. Statist.* 35 (6), 2313–2351.
- Chen, J., Li, H., 2013. Variable selection for sparse dirichlet-multinomial regression with an application to microbiome data analysis. *Ann. Appl. Stat.* 7 (1), 418–442.
- Efron, B., Hastie, T., Johnstone, I., Tibshirani, R., 2004. Least angle regression. *Ann. Statist.* 407–451.
- Fahrmeir, L., Kaufmann, H., 1985. Consistency and asymptotic normality of the maximum likelihood estimator in generalized linear models. *Ann. Statist.* 13 (1), 342–368.
- Fan, J., Li, R., 2001. Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Amer. Statist. Assoc.* 96 (456), 1348–1360.
- Friedman, J., Hastie, T., Tibshirani, R., 2010. Regularization paths for generalized linear models via coordinate descent. *J. Stat. Softw.* 33 (1), 1–22.
- Hastie, T., Tibshirani, R., Friedman, J., 2009. *The Elements of Statistical Learning: Data Mining, Inference and Prediction*. Springer, New York.
- Hoerl, A., Kennard, R., 1970. Ridge regression: biased estimation for nonorthogonal problems. *Technometrics* 12, 55–67.
- Krishnapuram, B., Carin, L., Figueiredo, M., Hartemink, A., 2005. Sparse multinomial logistic regression: fast algorithms and generalization bounds. *IEEE Trans. Pattern Anal. Mach. Intell.* 27, 957–968.
- McFadden, D., 1973. Conditional logit analysis of qualitative choice behaviour. In: Zarembka, P. (Ed.), *Frontiers in Econometrics*. Academic Press, New York.
- Meier, L., van de Geer, S., Bühlmann, P., 2008. The group lasso for logistic regression. *J. R. Stat. Soc. Ser. B* 70 (1), 53–71.
- Nemirovski, A., 1994. *Efficient methods in convex programming*. Lect. Notes.
- Nesterov, Y., 2004. *Introductory Lectures on Convex Optimization: A Basic Course*. Kluwer Academic Publishers.
- R Development Core Team, 2014. *R: a language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Simon, N., Friedman, J., Hastie, T., 2013. A blockwise descent algorithm for group-penalized multiresponse and multinomial regression. *arXiv Preprint*.
- Thurner, P., Eymann, A., 2000. Policy-specific alienation and indifference in the calculus of voting: a simultaneous model of party choice and abstention. *Public Choice* 102, 49–75.
- Tibshirani, R., 1996. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Ser. B* 58 (1), 267–288.
- Tibshirani, R., Saunders, M., Rosset, S., Zhu, J., Knight, K., 2005. Sparsity and smoothness via the fused lasso. *J. R. Stat. Soc. Ser. B* 67, 91–108.
- Turlach, B., Venables, W., Wright, S., 2005. Simultaneous variable selection. *Technometrics* 47, 349–363.
- Tutz, G., 2012. *Regression for Categorical Data*. Cambridge University Press, Cambridge.
- Tutz, G., Binder, H., 2006. Generalized additive modeling with implicit variable selection by likelihood-based boosting. *Biometrics* 62 (4), 961–971.
- Tvedebrink, T., 2010. Overdispersion in allelic counts and theta-correction in forensic genetics. *Theor. Popul. Biol.* 78, 200–210.
- Vincent, M., Hansen, N., 2014. Sparse group lasso and high dimensional multinomial classification. *Comput. Statist. Data Anal.* 71, 771–786.
- Wang, H., Leng, C., 2008. A note on adaptive group lasso. *Comput. Statist. Data Anal.* 52 (12), 5277–5286.
- Yellott, J., 1977. The relationship between Luce's choice axiom, Thurstone's theory of comparative judgement, and the double exponential distribution. *J. Math. Psych.* 15, 109–144.
- Yuan, M., Lin, Y., 2006. Model selection and estimation in regression with grouped variables. *J. R. Stat. Soc. Ser. B* 68 (1), 49–67.
- Zou, H., 2006. The adaptive lasso and its oracle properties. *J. Amer. Statist. Assoc.* 101 (476), 1418–1429.
- Zou, H., Hastie, T., 2005. Regularization and variable selection via the elastic net. *J. R. Stat. Soc. Ser. B* 67 (2), 301–320.
- Zou, H., Hastie, T., Tibshirani, R., 2007. On the “degrees of freedom” of the lasso. *Ann. Statist.* 35 (5), 2173–2192.