

Learning Transformation Models for Ranking and Survival Analysis

Vanya Van Belle

*Katholieke Universiteit Leuven, ESAT-SCD
Kasteelpark Arenberg 10
B-3001 Leuven, Belgium*

VANYA.VANBELLE@ESAT.KULEUVEN.BE

Kristiaan Pelckmans

*Uppsala University
Department of Information Technology
SysCon Polacksbacken
SE-751 05 Uppsala, Sweden*

KRISTIAAN.PELCKMANS@IT.UU.SE

Johan A. K. Suykens

Sabine Van Huffel

*Katholieke Universiteit Leuven, ESAT-SCD
Kasteelpark Arenberg 10
B-3001 Leuven, Belgium*

JOHAN.SUYKENS@ESAT.KULEUVEN.BE

SABINE.VANHUFFEL@ESAT.KULEUVEN.BE

Editor: Nicolas Vayatis

Abstract

This paper studies the task of learning transformation models for ranking problems, ordinal regression and survival analysis. The present contribution describes a machine learning approach termed MINLIP. The key insight is to relate ranking criteria as the Area Under the Curve to monotone transformation functions. Consequently, the notion of a Lipschitz smoothness constant is found to be useful for complexity control for learning transformation models, much in a similar vein as the 'margin' is for Support Vector Machines for classification. The use of this model structure in the context of high dimensional data, as well as for estimating non-linear, and additive models based on primal-dual kernel machines, and for sparse models is indicated. Given n observations, the present method solves a quadratic program existing of $O(n)$ constraints and $O(n)$ unknowns, where most existing risk minimization approaches to ranking problems typically result in algorithms with $O(n^2)$ constraints or unknowns. We specify the MINLIP method for three different cases: the first one concerns the preference learning problem. Secondly it is specified how to adapt the method to ordinal regression with a finite set of ordered outcomes. Finally, it is shown how the method can be used in the context of survival analysis where one models failure times, typically subject to censoring. The current approach is found to be particularly useful in this context as it can handle, in contrast with the standard statistical model for analyzing survival data, all types of censoring in a straightforward way, and because of the explicit relation with the Proportional Hazard and Accelerated Failure Time models. The advantage of the current method is illustrated on different benchmark data sets, as well as for estimating a model for cancer survival based on different micro-array and clinical data sets.

Keywords: support vector machines, preference learning, ranking models, ordinal regression, survival analysis

1. Introduction

Methods based on ranking continue to challenge researchers in different scientific areas, see, for example, Cl emen on et al. (2005), Herbrich, Graepel, and Obermayer (2000) and the references therein. Learning ranking functions offers a solution to different types of problems including ordinal regression, bipartite ranking and discounted cumulative gain ranking (DCG, see Cl emen on and Vayatis, 2007), studied frequently in research on information retrieval. These cases distinguish themselves in the definition (of the cardinality k) of the output domain and the chosen loss function. This paper deals with the general problem where the output domain can be arbitrary (with possibly infinite members $k = \infty$), but possesses a natural ordering relation between the members. Examples in which $k = \infty$ are found in survival analysis and preference learning in cases where the number of classes is not known in advance.

Earlier approaches to learning preference functions reduce the ranking problem to pairwise classification problems. This reasoning was followed in Ailon and Mohri (2008) and F urnkranz and H ullermeier (2003), Herbrich et al. (1998) and references therein. However, functions having high pairwise margins might still be bad approximations to real ranking problems. This is certainly the case in the (general) preference learning problem where possibly $k = \infty$: here a nonzero pairwise margin would need unnecessarily large parameters of the model. In this paper we address this issue by presenting a conceptual different approach: we adopt a smoothness condition on the ranking function to structure the space of ranking functions, and claim that this structure aligns in many applications better with the learning problems. This reasoning is motivated from relating a pairwise ranking criterion to a monotone *transformation* function. Besides empirical validation of this claim, we present formal relationships to other (statistical) models used for such tasks.

Figure 1 summarizes the ideas exposed in this work. First we describe the class of transformation models which contains two different components. The first component of a transformation model consists of a function $u : \mathbb{R}^d \rightarrow \mathbb{R}$ mapping the covariates $X \in \mathbb{R}^d$ to a value in \mathbb{R} such that the natural order on \mathbb{R} induces the ranking (approximately). Different names for such a function are found in literature depending on the problem setting, including a scoring, ranking, utility or health function. In this paper we will refer to this as to the utility function. The second component of the model maps this utility to an outcome in \mathbb{R} by a transformation function $h : \mathbb{R} \rightarrow \mathbb{R}$. This is a univariate monotonically increasing function, basically capturing the scale of the output. The central observation now is that when one knows the ordinal relations between instances, one can estimate a transformation function mapping the instances to their utility value $\hat{u}(X)$. Depending on the problem at hand one is interested in the results of the first or second component of the transformation model. For ranking and survival analysis one typically ignores the second phase, whereas in ordinal regression a prediction of the output level is found by combining the first and the second components.

Transformation models are especially appropriate when considering data arising from a survival study. Survival analysis concerns data which represent a time-to-event, as, for example, a patient relapsing after surgery, or the time till a part of a mechanical device breaks down, see Kalbfleisch and Prentice (2002) for a broad survey of this field. The goal in survival analysis is often to relate time-to-event of an instance to a corresponding set of covariates. While practice and theoretical results here continue to have a strong impact in most quantitative scientific areas, survival analysis has been studied only sporadically in a context of machine learning, and such studies are mostly found in the field of artificial neural networks, see, for example, Biganzoli et al. (1998) and Kattan

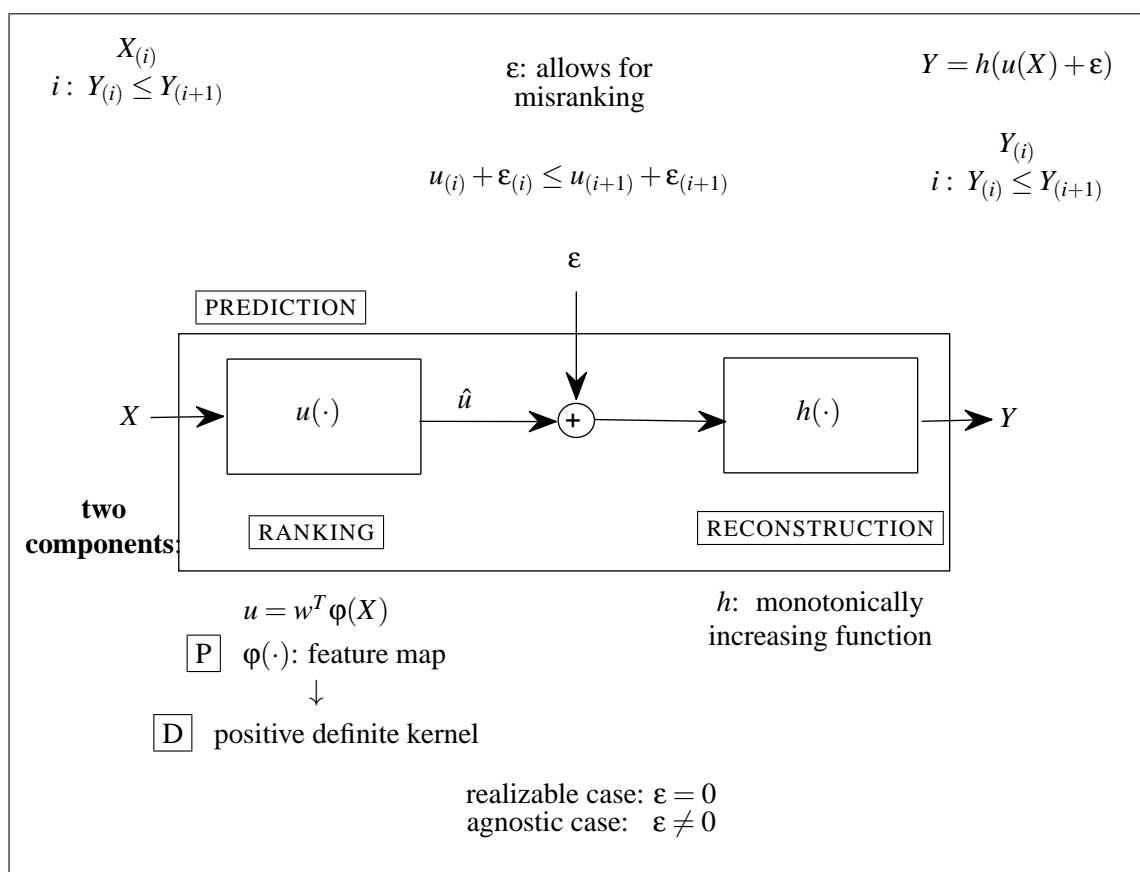


Figure 1: Overview: Transformation models consist of two components, the *utility function* u and a *transformation function* h . Given a data set $\mathcal{D} = \{(X_{(i)}, Y_{(i)})\}_{i=1}^n$ where the instances are sorted such that $Y_{(i)} \leq Y_{(i+1)}$, a utility function $u(X) = w^T \varphi(X)$ is trained such that the ranking on the evaluations of this function is representative for the ranking on the outcome. In the *realizable case* the ordering in utility will exactly coincide with the ordering in observed outcome $\{Y_{(i)}\}_i$. In the *agnostic case* however, the ordering will only be exact up to appropriate (nonzero) error variables $\{\varepsilon_i\}_i$. The modelling procedure will also be performed in two steps. The first step recovers u ('ranking'), while the second step is concerned with learning an explicit representation of the transformation function ('reconstruction'). In practice (depending on the problem at hand) one is mostly interested in implementing the first step only.

et al. (1997). However, we claim that there is a large potential for such studies: as (i) the approach of classical likelihood-based approaches have their intrinsic limitations, especially when a realistic underlying model cannot be assumed. A distribution-free approach is more appropriate here; (ii) A risk-based approach is often easier as one does not care about recovering the exact parameters describing the process of interest, but one is only interested in making good predictions, or exploring structure in the data; (iii) Computational issues for the classical statistical approach persist, and the

question how to solve the estimation equations numerically is often approached in an ad hoc way (if at all, see Kalbfleisch and Prentice, 2002 and references).

We find that the class of transformation models is a powerful tool to model data arising from survival studies for different reasons. The first reason being that they separate nicely the model for the time-scale (via the transformation function), and the qualitative characterization of an instance (via the utility function). We will furthermore argue on the close relationship with existing techniques as Cox' proportional hazard and accelerated failure time (AFT) models, see, for example, Dabrowska and Doksum, 1988, Koenker and Geling, 2001, Cheng, Wei, and Ying, 1997 and citations. In the following, we will relate the transformation function to ranking criteria as Kendall's τ or area under the curve (AUC), hence outlining a unified framework to study survival models as used in a statistical context and machine learning techniques for learning ranking functions. This relation indicates how one may apply the method of structural risk minimization (SRM, see Vapnik, 1998) here. The immediate consequence is the possibility to apply learning theory, with the capabilities to explain good performances in modelling high-dimensional data sets as well as for non-linear models (see Vapnik, 1998). Thirdly, in studies of failure time data, *censoring* is omnipresent. Censoring prohibits that one observes the actual event of interest fully, but gives partial information on the outcome instead. The prototypical case is 'a patient hasn't suffered the event as yet, but may experience an event in the future', but many other examples are studied. We will see how the proposed approach can handle censored observations conveniently.

The computational merit of this paper is then how one can fit such a model efficiently to data. Therefore, we consider an appropriate class of utility functions, either linear functions, or kernel based models. Secondly, instead of restricting attention to a parameterized class of transformation functions, we let the transformation function of interest be unspecified as one does for partial likelihood approaches, see Kalbfleisch and Prentice (2002). Especially, we define the appropriate transformation function only on the observed samples, by inferring an appropriate set of ordinal relations between them. Then we observe that the Lipschitz smoothness constant associated to such a transformation function can also be evaluated based on the samples only. Consequently, our fitting strategy called MINLIP finds the maximally smooth (implicitly defined) transformation function fitting the data samples. This is the *realizable case* where we can make the assumption of the existence of such a transformation model. In case we allow for misfit, we extend the model using slack-variables. It is then found that this problem can be solved as a convex Quadratic Program (QP), for which highly efficient software is readily available. In the case of utility functions which are kernel based models, we indicate how one can represent the solution as a sum of positive definite kernels, and the Lagrange dual problem again solves the corresponding problem as a convex QP. For the case linear utility functions are considered, we suggest how one can obtain zero parameters ('sparseness') suggesting structure in the data using an 1-norm regularization mechanism (see also Tibshirani, 1996).

Besides the conceptual and computational discussion, this paper gives empirical evidence for the approach. We consider empirical studies of ordinal regression and survival analysis. Performance of MINLIP on ordinal regression is analyzed using the ordinal data compiled by Chu and Keerthi (2005). MINLIP is applied on two different survival studies. A first study involves micro-array data sets: two breast cancer data sets (Sørliet et al., 2003; van Houwelingen et al., 2006) and one data set concerning diffuse large B-cell carcinoma (Rosenwald et al., 2002). In a last study, concerning a clinical breast cancer survival study (Schumacher et al., 1994), we investigate the estimation of

non-linear covariate effects and compare results obtained with MINLIP with Cox regression with penalized smoothing splines.

In Van Belle et al. (2007) we proposed a modification to standard SVMs to handle censored data. A computationally less demanding algorithm was presented in Van Belle et al. (2008). Starting from this latter model, we replaced the maximal margin strategy with the minimal Lipschitz smoothness strategy as presented in Van Belle et al. (2009). This work extends considerably the results of this short paper. Most notably, this paper additionally elaborates on the case of survival analysis and a number of new case studies. The different application areas in which the proposed method can be applied are summarized in Table 1. In addition, it is stated how the model needs to be used and which equations need to be solved to obtain the solution.

This paper is organized as follows. The following Section discusses in some detail the use of transformation models and its relation with ranking methods. Section 3 studies the estimator in a context of ranking. Section 4 specifies how MINLIP is to be used in a context of ordinal regression, where only k different output levels are possible. Section 5 discusses the use of MINLIP in the presence of censoring. In Section 6 experiments illustrate the use of the MINLIP method.

2. Transformation Models and Ranking Methods

In this paper we work in a stochastic context, so we denote random variables as capital letters, for example, X, Y, \dots , which follow an appropriate stochastic law P_X, P_Y, \dots , abbreviated (generically) as P . Deterministic quantities as constants and functions are represented in lower case letters (e.g., d, h, u, \dots). Matrices are denoted as boldface capital letters (e.g., $\mathbf{X}, \mathbf{D}, \dots$). Ordered sets will be denoted as $\{S_{(i)}\}$, indicating that $S_{(i)} \leq S_{(i+1)}$. Before the relation between transformation models and ranking methods can be explored, some terminology needs to be defined.

Definition 1 (Lipschitz smoothness) *A univariate function $h(Z)$ has a Lipschitz constant $L \geq 0$ if*

$$|h(Z) - h(Z')| \leq L|Z - Z'|, \forall Z, Z' \in \mathbb{R}.$$

A *transformation model* is then defined as follows:

Definition 2 (Transformation Model) *Let $h : \mathbb{R} \rightarrow \mathbb{R}$ be a strictly increasing function with Lipschitz constant $L < \infty$, and let $u : \mathbb{R}^d \rightarrow \mathbb{R}$ be a function of the covariates $X \in \mathbb{R}^d$. Let ε be a random variable ('noise') independent of X , with cumulative distribution function $F_\varepsilon(e) = P(\varepsilon \leq e)$ for any $e \in \mathbb{R}$. Then a Noisy Transformation Model (NTM) takes the form*

$$Y = h(u(X) + \varepsilon). \tag{1}$$

In the remainder of the paper, we will use Z to denote $u(X) + \varepsilon$ for notational convenience. Now the problem is reduced to estimating a utility function $u : \mathbb{R}^d \rightarrow \mathbb{R}$ and a transformation function h from a set of i.i.d. observations $\{(X_i, Y_i)\}_{i=1}^n$ without imposing any distributional (parametric) assumptions on the noise terms $\{\varepsilon_i\}$. Note that without structural assumptions, the utility can not uniquely be defined. Later on, we will specify similar assumptions as in the maximal margin strategy of Vapnik when introducing support vector machines, to find a unique solution for the utility function.

Task	Subtasks	algorithm	necessary training data for algorithm	result of algorithm	equation	comment
ranking	ranking	MINLIP	$\{(X_i, Y_i)\}$	$\hat{u}(X)$	(8)	Use in combination with (10)* to obtain sparse models
ordinal regression	ranking	MINLIP	$\{(X_i, Y_i)\}, \mathcal{B}$	$\hat{u}(X), \hat{\nu}$	(11)	Use in combination with (10)* to obtain sparse models
	reconstruction	comparison with thresholds ν	$\{\hat{u}(X_i)\}, \hat{\nu}$	\hat{Y}	Figure 5	The goal is to predict class label
survival analysis	ranking	MINLIP	$\{(X_i, Y_i)\}$	$\hat{u}(X)$	(16)	Use in combination with (10)* to obtain sparse models
	reconstruction	monotonic regression after replication in consecutive time intervals	$\{(\hat{u}(X_i), Y_i)\}$	$\{\hat{Y}_{it}\}$		The goal is to predict hazard and/or survival function

*Equation (10) can only be used in combination with a linear kernel

Table 1: Overview of methods and applications proposed in the paper. Depending on the problem at hand, a different version of the proposed model needs to be applied. Depending on the subtasks, different training data (variables X_i , target value Y_i , utility $u(X_i)$), dummy responses \mathcal{B} , ... need to be given to the algorithm to obtain the desired response (utility $u(X_i)$, transformation function $h(u(X))$), prediction $\hat{Y}_i = h(u(X_i))$, threshold values $\hat{\nu}$, risk on event within the l^{th} interval \hat{Y}_{it}, \dots .

Kalbfleisch and Prentice (2002) considered transformation models for failure time models. The transformation models discussed in Cheng, Wei, and Ying (1997), Dabrowska and Doksum (1988) and Koenker and Geling (2001) differ from the above definition in the transformation function h . They define the model as $h^{-}(Y) = u(X) + \varepsilon$, which is equivalent to (1) if $h^{-}(h(Z)) = h(h^{-}(Z)) = Z$ for all Z .

To relate transformation models with ranking functions, we reason as follows. To express the performance of a ranking function one can use Kendall’s τ , area under the curve (AUC) or a related measure. In this paper we will work with the *concordance* of a function $u : \mathbb{R}^d \rightarrow \mathbb{R}$ respective to the outcome. The concordance is defined as the probability that the order in outcome of two i.i.d. observations (X, Y) and (X', Y') is preserved in the utility u :

$$C(u) = P((u(X) - u(X'))(Y - Y') > 0). \tag{2}$$

Given a set of n i.i.d. observations $\{(X_i, Y_i)\}_{i=1}^n$, the empirical concordance index is then calculated as

$$C_n(u) = \frac{2}{n(n-1)} \sum_{i < j} I[(u(X_i) - u(X_j))(Y_i - Y_j) > 0],$$

where the indicator function $I(z)$ equals 1 if $z > 0$, and equals zero otherwise. Equivalently, the risk is defined as follows.

Definition 3 (Risk of (h, u)) *The risk associated with a monotonically increasing function penalizes discordant samples $h(u(X))$ and $h(u(X'))$ as*

$$\mathcal{R}(u) = P((h(u(X)) - h(u(X')))(Y - Y') < 0).$$

Or, since h is monotonically increasing, the risk is expressed as

$$\mathcal{R}(u) = P((u(X) - u(X'))(Y - Y') < 0).$$

Its empirical counterpart then becomes

$$\mathcal{R}_n(u) = 1 - C_n(u).$$

Empirical Risk Minimization (ERM) is then performed by solving

$$\hat{u} = \arg \min_{u \in \mathcal{U}} \mathcal{R}_n(u) = \arg \max_{u \in \mathcal{U}} C_n(u), \tag{3}$$

where $\mathcal{U} \subset \{u : \mathbb{R}^d \rightarrow \mathbb{R}\}$ is an appropriate subset of ranking functions, see, for example, Cl emen on et al. (2005) and citations. However, this approach results in combinatorial optimization problems. One therefore majorizes the discontinuous indicator function by the Hinge loss, that is, $\ell(z) \leq \max(0, 1 - z)$ yielding rankSVM (Herbrich, Graepel, and Obermayer, 2000). The disadvantage of this solution is that it leads to $O(n^2)$ number of constraints or unknowns, often making it difficult to apply to real life problems. A solution to this problem is found in relating transformation models with Equation (3): if a function $u : \mathbb{R}^d \rightarrow \mathbb{R}$ exists such that $C_n(u) = 1$, one describes implicitly a transformation function (see Figure 2). If two variables u and y are perfectly concordant, then there exists a monotonically increasing function h such that $h(u)$ and y are perfectly concordant. Moreover, there exists such a function h , with Lipschitz constant L , mapping u to y such that $y = h(u)$. Or more formally:

Lemma 1 (Existence of a Transformation Function) *Given a collection of pairs $\{(Z_{(i)}, Y_{(i)})\}_{i=1}^n$, enumerated such that $Y_{(i)} \leq Y_{(j)}$ if and only if $i \leq j$, and considering the conditions on the observations for $L < \infty$:*

$$0 \leq Y_{(i)} - Y_{(j)} \leq L(Z_{(i)} - Z_{(j)}), \quad \forall i < j = 1, \dots, n, \quad (4)$$

we state that:

1. *If one has for a finite value $L \geq 0$ that (4) holds, then there exists a monotonically increasing function $h : \mathbb{R} \rightarrow \mathbb{R}$ with Lipschitz constant L interpolating the data points.*
2. *If for all admissible $(Z, Y) \in \mathbb{R} \times \mathbb{R}$ one has that $Y = h(Z)$ for an (unknown) continuous, (finite) differentiable and monotonically increasing function $h : \mathbb{R} \rightarrow \mathbb{R}$, then there is a value $L < \infty$ such that (4) holds.*

Proof To prove 1, consider the linear interpolation function $h_n : \mathbb{R} \rightarrow \mathbb{R}$, defined as

$$h_n(Z) = \frac{Z - Z_{\bar{z}(Z)}}{Z_{\bar{z}(Z)} - Z_{\underline{z}(Z)}} (Y_{\bar{z}(Z)} - Y_{\underline{z}(Z)}) + Y_{\underline{z}(Z)},$$

where we define $\bar{z}(Z) = \arg \min_{i \in \{1, \dots, n\}} (Z_i : Z_i > Z)$ and $\underline{z}(Z) = \arg \max_{i \in \{1, \dots, n\}} (Z_i : Z_i \leq Z)$. Direct manipulation shows that this function is monotonically increasing and continuous. Now take $Z < Z' \in \mathbb{R}$, then we have to show that $h_n(Z') - h_n(Z) \leq L(Z' - Z)$. For notational convenience define $l = \underline{z}(Z)$, $u = \bar{z}(Z)$, $l' = \underline{z}(Z')$ and $u' = \bar{z}(Z')$, then

$$\begin{aligned} h_n(Z') - h_n(Z) &= \frac{Z' - Z_{l'}}{Z_{u'} - Z_{l'}} (Y_{u'} - Y_{l'}) + Y_{l'} - \frac{Z - Z_l}{Z_u - Z_l} (Y_u - Y_l) - Y_l \\ &\leq L(Z' - Z_{l'}) - L(Z - Z_l) + L(Z_{l'} - Z_l) \\ &= L(Z' - Z), \end{aligned}$$

where we use that $Y_{l'} - Y_l \leq L(Z_{l'} - Z_l)$.

Item 2 is proven as follows. Let such an h exist, then the mean value theorem asserts that for any two samples (Z_i, Y_i) and (Z_j, Y_j) for which $Z_i \leq Z_j$, there exists a Z within the interval $[Z_i, Z_j] \subset \mathbb{R}$ such that

$$(Y_j - Y_i) = (Z_j - Z_i)h'(Z) \leq L(Z_j - Z_i),$$

where $L = \sup_Z h'(Z)$.

Note that Equation (4) implies that $C_n(Z) = 1$. ■

3. MINLIP: A Convex Approach to Learning a Transformation Model

This Section describes how transformation models can be learned by means of a convex approach. The Section starts with a discussion of the realizable case and extends this model formulation towards the agnostic case and non-linearities using Mercer kernels.

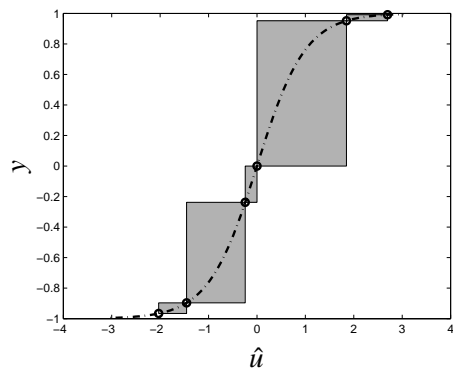


Figure 2: Relation between ranking and transformation models: if two variables u and y are perfectly concordant, they describe a monotonically increasing function $y = h(u)$. The dots represent u and outcome y for training points. In these observations the value of the function h is known exactly. To predict the y -value of the test observations, the function h needs to be approximated between the training points (grey area). All functions \hat{h} which are monotonically increasing and lie within the grey zones are valid prediction rules.

3.1 The Realizable Case

The realizable case refers to the situation where there exists a function $u(X)$ such that the ranking of $u(X)$ perfectly reflects the ranking of Y . Otherwise stated, there exists a function $u(X)$ such that $C(u) = 1$. Lemma 1 describes the existence of h , but since this transformation function is only known at the training points, it is not unique. Figure 2 illustrates that all monotonically increasing functions lying within the grey bounds satisfy the conditions. Therefore, the Lipschitz constant is used to control the complexity of the transformation function. Transformation functions with a smaller Lipschitz constant will be preferred. For notational convenience, we will assume no coinciding outcomes (ties). Let h be a monotonically increasing function with Lipschitz constant $L < \infty$, such that $h(Z) - h(Z') \leq L(Z - Z')$ for all $Z \geq Z'$. Restricting attention to the observations $\{(X_{(i)}, Y_{(i)})\}_{i=1}^n$, one has the necessary conditions

$$h(u(X_{(i)})) - h(u(X_{(i-1)})) \leq L(u(X_{(i)}) - u(X_{(i-1)})) ,$$

for all $i = 2, \dots, n$. Here, we assume that the data obey a noiseless transformation model ($\epsilon = 0$ in (1)). For now, linear utility functions defined as

$$u(X) = w^T X ,$$

are considered. Extensions towards non-linear utility functions using Mercer kernels are handled in Subsection 3.3. Since the function $u(X) = w^T X$ can be arbitrary rescaled such that the corresponding transformation function has an arbitrary Lipschitz constant (i.e., for any $\alpha > 0$, one has $h(u(X)) = \tilde{h}(\tilde{u}(X))$ where $\tilde{h}(Z) \triangleq h(\alpha^{-1}Z)$ and $\tilde{u}(X) = \alpha u(X)$), we fix the norm $w^T w$ and try to find $u(X) = v^T X$ with $v^T v = 1$. Hence learning a transformation model with minimal Lipschitz

constant of h can be written as

$$\begin{aligned} \min_{v,L} \quad & \frac{1}{2}L^2 \\ \text{s.t.} \quad & \begin{cases} \|v\|_2 = 1 \\ Y_{(i)} - Y_{(i-1)} \leq L(v^T X_{(i)} - v^T X_{(i-1)}), \quad \forall i = 2, \dots, n. \end{cases} \end{aligned}$$

Substituting $w = Lv$ we get equivalently:

$$\begin{aligned} \min_w \quad & \frac{1}{2}w^T w \\ \text{s.t.} \quad & Y_{(i)} - Y_{(i-1)} \leq w^T X_{(i)} - w^T X_{(i-1)}, \quad \forall i = 2, \dots, n, \end{aligned} \tag{5}$$

which goes along similar lines as the hard margin SVM (see, e.g., Shawe-Taylor and Cristianini, 2004) and ranking SVM (Freund et al., 2004), where the threshold value 1 is replaced by $Y_{(i)} - Y_{(i-1)}$. Note that an intercept term is not needed since differences in utility are used. Observe that this problem has $n - 1$ linear constraints. We will refer to this model as MINLIP.

Problem (5) can be compactly rewritten as

$$\begin{aligned} \min_w \quad & \frac{1}{2}w^T w \\ \text{s.t.} \quad & \mathbf{D}\mathbf{X}w \geq \mathbf{D}\mathbf{Y}, \end{aligned}$$

where $\mathbf{X} \in \mathbb{R}^{n \times d}$ is a matrix with each row containing one observation, that is, $\mathbf{X}_i = X_{(i)} \in \mathbb{R}^d$ and $\mathbf{Y} = [Y_{(1)} \cdots Y_{(n)}]^T$, a vector with the corresponding outcomes. The matrix $\mathbf{D} \in \{-1, 0, 1\}^{(n-1) \times n}$

$$\mathbf{D} = \begin{bmatrix} -1 & 1 & 0 & 0 & \dots & 0 & 0 \\ 0 & -1 & 1 & 0 & \dots & 0 & 0 \\ \vdots & & & & & & \vdots \\ 0 & \dots & 0 & 0 & 0 & -1 & 1 \end{bmatrix},$$

gives the first order differences of a vector, that is, assuming no ties in the output, $\mathbf{D}_i \mathbf{Y} = Y_{(i+1)} - Y_{(i)}$ for all $i = 1, \dots, n - 1$, with \mathbf{D}_i the i^{th} row of \mathbf{D} .

In the presence of ties, $Y_{(i+1)}$ is replaced by $Y_{(j)}$, with j the smallest output value with $Y_{(j)} > Y_{(i)}$. See Section 4 for more details. Solving this problem as a convex QP can be done efficiently with standard mathematical solvers as implemented in MOSEK¹ or R-quadprog.² The following proposition states when the MINLIP model is valid.

Proposition 1 (Validity of MINLIP) *Assume that $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ would obey the relation*

$$Y = h_0(w_0^T X), \tag{6}$$

where we refer to the (fixed but unknown) vector $w_0 \in \mathbb{R}^d$ as to the 'true' parameters, and to the (fixed but unknown) monotonically increasing function $h_0 : \mathbb{R} \rightarrow \mathbb{R}$ as the 'true' transformation function. Let for each couple (X, Y) and (X', Y') where $Y \neq Y'$ the constant $L' > 0$ be defined as

$$\frac{1}{L'} = \frac{w_0^T (X - X')}{Y - Y'},$$

1. MOSEK can be found at <http://www.mosek.org>.

2. R-quadprog can be found at <http://cran.r-project.org/web/packages/quadprog/index.html>.

where $L' = \infty$ if $w_0^T(X - X') = 0$. By construction we have that $L' \leq L_0$ and that the constant exists everywhere. The result of the MINLIP model then becomes:

$$\mathcal{L} = \max_{\|w\|_2=1} \min_{Y > Y'} \frac{w^T(X - X')}{Y - Y'} = \max_{\|w\|_2=1} \min_{Y > Y'} \frac{w^T(X - X')}{L' w_0^T(X - X')}. \quad (7)$$

We then state that MINLIP yields a good approximation of the parameter vector w_0 in the noiseless case as long as there are enough observations (X, Y) such that $w_0^T(X - X') \approx 1$ and $L' \approx L_0$.

Proof Let the unit-length vector $\overline{(X - X')} \in \mathbb{R}^d$ be defined as $X - X' = \overline{(X - X')} \|X - X'\|_2$, then we can write (7) as

$$\max_{\|w\|_2=1} \min_{Y > Y'} \frac{w^T \overline{(X - X')}}{L' w_0^T \overline{(X - X')}}.$$

Let us now focus attention on the set $\mathcal{S} = \{(X - X', Y - Y') : (X, Y), (X', Y') \in \mathcal{D} = \{X_i, Y_i\}_{i=1}^n\}$, where $\mathcal{L} = w^T(X - X')/(Y - Y')$ for which this value \mathcal{L} is actually achieved. It is seen that the estimate w lies in the span of this set \mathcal{S} as otherwise the maximum value could be increased. When we assume that the data set contains enough observations (X, Y) such that $w_0^T(X - X') \approx 1$ and $L' \approx L_0$, they will end up in the set \mathcal{S} , and as a result we have that $w^T w_0 \approx 1$. As the optimal solution is fully determined by the terms $w_0^T \overline{(X - X')} \approx 1$ and $L' \approx L_0$ (cfr. duality results in convex optimization), one should also have that $w \approx w_0$. ■

Formally, consistency of MINLIP in the asymptotic case under a sufficient condition of the data being non-degenerate is derived in Appendix A.

3.2 The Agnostic Case

In case it is impossible to find a utility function $u : \mathbb{R}^d \rightarrow \mathbb{R}$ extracting the ranking perfectly, a noisy transformation model is considered:

$$Y = h(w^T X + \varepsilon),$$

where $u = w^T X$. The introduction of the error variable asks for an adaptation of the Lipschitz-based complexity control. As a loss function $\ell : \mathbb{R} \rightarrow \mathbb{R}$ we choose the absolute value loss $\ell(\varepsilon) = |\varepsilon|$ for three reasons: (i) It is known that this loss function is more robust to misspecification of the model and outliers than, for example, the squared loss $\ell(\varepsilon) = \varepsilon^2$; (ii) The use of the absolute value loss will result in sparse solutions with many error terms equal to zero; (iii) In binary classification this norm is well performing in SVMs. However, the choice of the loss remains arbitrary. Incorporation of the errors (slack variables) leads to the following model formulation:

$$\begin{aligned} \min_{w, \varepsilon} \quad & \frac{1}{2} w^T w + \gamma \|\varepsilon\|_1 \\ \text{s.t.} \quad & \mathbf{D}(\mathbf{X}w + \varepsilon) \geq \mathbf{D}\mathbf{Y}, \end{aligned} \quad (8)$$

where $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T \in \mathbb{R}^n$ represents the errors, $\|\varepsilon\|_1 = \sum_{i=1}^n |\varepsilon_i|$ and $\gamma > 0$ is a regularization constant, making a trade-off between model complexity and error. This problem can again be solved as a convex quadratic program.

3.3 A Non-linear Extension using Mercer Kernels

Let $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}^{d_\varphi}$ be a feature map mapping the data to a high dimensional feature space (of dimension d_φ , possibly infinite). A non-linear utility function can then be defined as

$$u(X) = w^T \varphi(X),$$

with $w \in \mathbb{R}^{d_\varphi}$ a vector of unknowns (possibly infinite dimensional). Take $\Phi = [\varphi(X_{(1)}), \dots, \varphi(X_{(n)})]^T \in \mathbb{R}^{n \times d_\varphi}$. The realizable learning problem can then be represented as:

$$\begin{aligned} \min_w \quad & \frac{1}{2} w^T w \\ \text{s.t.} \quad & \mathbf{D}\Phi w \geq \mathbf{D}\mathbf{Y}, \end{aligned}$$

with the matrix \mathbf{D} defined as before. The Lagrange dual problem becomes

$$\begin{aligned} \min_\alpha \quad & \frac{1}{2} \alpha^T \mathbf{D}\mathbf{K}\mathbf{D}^T \alpha - \alpha^T \mathbf{D}\mathbf{Y} \\ \text{s.t.} \quad & \alpha \geq \mathbf{0}_{n-1}, \end{aligned}$$

where the kernel matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$ contains the kernel evaluations such that $\mathbf{K}_{ij} = \varphi(X_i)^T \varphi(X_j)$ for all $i, j = 1, \dots, n$. The estimated utility \hat{u} can be evaluated at any point $X^* \in \mathbb{R}^d$ as

$$\hat{u}(X^*) = \hat{\alpha}^T \mathbf{D}\mathbf{K}_n(X^*), \quad (9)$$

where $\mathbf{K}_n(X^*) = [K(X_1, X^*), \dots, K(X_n, X^*)]^T \in \mathbb{R}^n$. The dual (Shawe-Taylor and Cristianini, 2004; Suykens, Gestel, Brabanter, Moor, and Vandewalle, 2002; Vapnik, 1998) of the agnostic learning machine of Subsection 3.2 is obtained analogously:

$$\begin{aligned} \min_\alpha \quad & \frac{1}{2} \alpha^T \mathbf{D}\mathbf{K}\mathbf{D}^T \alpha - \alpha^T \mathbf{D}\mathbf{Y} \\ \text{s.t.} \quad & \begin{cases} -\gamma \mathbf{1}_n \leq \mathbf{D}^T \alpha \leq \gamma \mathbf{1}_n \\ \alpha \geq \mathbf{0}_{n-1}, \end{cases} \end{aligned}$$

with \mathbf{K} as above and the resulting estimate can be evaluated as in (9) without computing explicitly \hat{w} nor φ . We refer to Appendix B for a detailed derivation. Typical choices for kernel functions are:

$$\begin{aligned} K(X, X_i) &= X_i^T X \quad (\text{linear kernel}) \\ K(X, X_i) &= (\tau + X_i^T X)^d, \quad \tau \geq 0 \quad (\text{polynomial kernel of degree } d) \\ K(X, X_i) &= \exp\left(-\frac{\|X - X_i\|_2^2}{\sigma^2}\right) \quad (\text{RBF kernel}). \end{aligned}$$

In cases where one is interested in the modelling of covariate effects, one could use an additive utility function:

$$u(X) = \sum_{p=1}^d u^p(X^p),$$

where X^p represents the p^{th} covariate of datapoint X . Using Equation (9) this can be written as:

$$\begin{aligned}\hat{u}(X) &= \sum_{p=1}^d \alpha^T \mathbf{D} \mathbf{K}^p(X^p) \\ &= \alpha^T \mathbf{D} \sum_{p=1}^d \mathbf{K}^p(X^p),\end{aligned}$$

where the kernel matrix $\mathbf{K}^p \in \mathbb{R}^{n \times n}$ contains the kernel evaluations such that $\mathbf{K}_{ij}^p = \phi(X_i^p)^T \phi(X_j^p)$ for all $i, j = 1, \dots, n$. As a result, componentwise kernels (Pelckmans et al., 2005b):

$$K(X, X_i) = \sum_{p=1}^d K^p(X^p, X_i^p),$$

which can be seen as a special case of ANOVA kernels (Vapnik, 1998), can be used. The use of such componentwise kernels allows for interpreting the non-linear effects of the covariates.

3.4 Prediction with Transformation Models

Prediction of the outcome using transformation models is a two-step approach (see Figure 1). First, the utility $u(X)$ is estimated, giving an ordering relation between the observations. When interested in an outcome prediction, the transformation function h has to be estimated. The prediction step is a univariate regression problem, which can be solved using monotonic regression models. Remark that in the ranking setting, one is not interested in the estimation of the transformation function since the goal is to find the ranking. Estimation of the transformation function for ordinal regression and survival analysis will be illustrated later.

3.5 Toward Sparse Solutions using $\|w\|_1$

This subsection describes an extensions to the above model. Specifically, we will be interested in the case where d is large compared to n . Consequently, we will be interested in computational methods which reveal the relevant input variables of use in the learned prediction rule. We restrict ourselves to the primal case where $u(X) = w^T X$ for the linear case and an unknown monotonically increasing function $h : \mathbb{R} \rightarrow \mathbb{R}$. In this extension an l_1 penalty (Tibshirani, 1996) is used instead of the term $w^T w$. We shall refer to this model as MINLIP $_{L1}$:

$$\begin{aligned}\min_{w, \varepsilon} \quad & \|w\|_1 + \gamma \|\varepsilon\|_1 \\ \text{s.t.} \quad & \mathbf{D}(\mathbf{X}w + \varepsilon) \geq \mathbf{D}\mathbf{Y},\end{aligned}\tag{10}$$

where $\|w\|_1 = \sum_{p=1}^d |w_p|$. This linear programming problem (LP) can be solved efficiently with standard mathematical solvers. This formulation does not allow for a straightforward dual derivation.

Figure 3 illustrates the possible advantage of the sparse alternative over the standard MINLIP formulation. We created 100 artificial data sets, each containing 150 observations with 200 covariates. 100 observations were used for training, the remaining for testing. A varying number of $d = 100, 110, \dots, 200$ covariates were used to build the outcome, all other features being irrelevant. All covariates were drawn from a normal distribution with zero mean and standard deviation 1. The outcome was obtained as a weighted sum of the relevant covariates, where the weights were drawn from a standard normal distribution. The test error of the MINLIP $_{L1}$ model was lower than for the standard model.

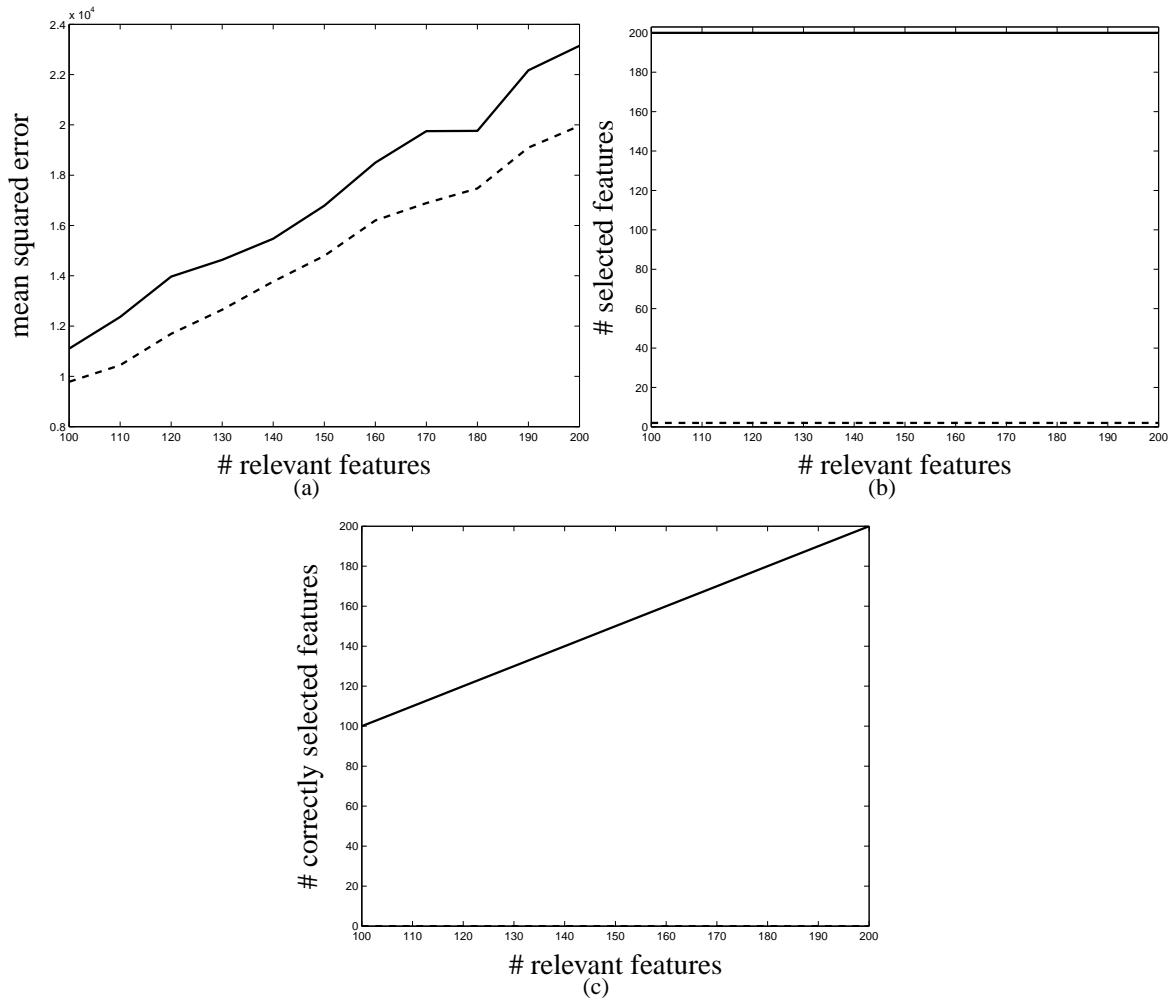


Figure 3: Performance and feature selection ability of MINLIP (solid) and MINLIP_{L1} (dashed) on an artificial data set ($n=100$ for training, $n_{\text{test}}=50$ for testing). 200 $\mathcal{N}(0, 1)$ distributed covariates were generated, a varying number $d = 100, 110, \dots, 200$ of which were used to generate the outcome ($Y = \sum_{p=1}^d w^p X^p$, with w drawn from a standard normal distribution). The results are averaged over 100 data sets. (a) Median mean squared error on the test sets: MINLIP_{L1} performs better than MINLIP. (b-c) Number of selected (absolute value of estimated weight $> 10^{-8}$) and correctly selected variables versus number of relevant variables. The MINLIP method selects all variables, a lot of them not being relevant. The MINLIP_{L1} model selects very few variables, but those which are selected are also relevant.

3.6 Comparison with Other Methods

An approach often seen within preference ranking problems is the reformulation of the ranking problem as a classification problem. Examples of this strategy can be found in Ailon and Mohri (2008), Fürnkranz and Hüllermeier (2003) and Herbrich et al. (1998). However, transforming ranking to classification deflects attention from the underlying problem within ranking problems. In contrast with these methods, the MINLIP approach concentrates on the ranking problem by use of the transformation model.

Currently used ranking methods include ranksVM (Herbrich, Graepel, and Obermayer, 2000) and RankBoost (Freund et al., 2004). Although the method proposed here and ranksVM are both based on SVMs, two differences can be noted: (i) firstly, the ranksVM uses all pairs of data points for training, which results in $O(n^2)$ comparisons, where MINLIP has a complexity of $O(n)$. This reduction in complexity makes the model more applicable to large data sets; (ii) Secondly, the complexity control, being the margin and the Lipschitz constant, is different in both methods. In ranksVM all margins are equal and the model is tuned to maximize this margin. In MINLIP the margins differ corresponding to the difference in the output levels.

4. Learning for Ordinal Regression

Consider now the situation where the output takes a finite number of values - say $k \in \mathbb{N}$ - and where the k different classes possess a natural ordering relation. In this case the outcome Y is an element of the finite ordered set $\{Y_{(1)}, \dots, Y_{(k)}\}$.

4.1 A Modification to MINLIP

In Section 3.1 it is mentioned that comparisons are made between points (i) and (j) where $Y_{(j)}$ is the first ordered value bigger than $Y_{(i)}$. Applying this methodology in the ordinal setting would lead to as many comparisons with point (i) from class k_i as there are observations in class $k_i + 1$. To cope with this issue, we add dummy observations (X, B) in between two consecutive ordinal classes with levels $Y_{(i)} < Y_{(i+1)}$ such that $B_{(i)} = \frac{1}{2}(Y_{(i+1)} + Y_{(i)})$ (see Figure 4) and leaving their covariates and utility function unspecified. This implies that one has to compare each observation only twice, once with the dummy observation in between the previous and the current ordinal class and once with the dummy observation in between the current and the next class, restricting the number of constraints to $O(n)$. The solution of this problem can be found implicitly by extending the $\mathbf{Y} \in \mathbb{R}^n$ and $\mathbf{X} \in \mathbb{R}^{n \times d}$ matrices as follows:

$$\tilde{\mathbf{Y}} = \begin{bmatrix} \mathbf{Y} \\ B_{(1)} \\ B_{(2)} \\ \dots \\ B_{(k-1)} \end{bmatrix} \quad \text{and} \quad \tilde{\mathbf{X}} = \left[\begin{array}{c|c} \mathbf{X} & 0 \\ \hline 0 & I_{k-1} \end{array} \right],$$

where $\tilde{\mathbf{X}} \in \mathbb{R}^{(n+k-1) \times (d+k-1)}$ and $\tilde{\mathbf{Y}} \in \mathbb{R}^{n+k-1}$ and I_{k-1} represents the identity matrix of dimension $k - 1$. The problem is then formulated as in Equation (8) after replacing \mathbf{X} by $\tilde{\mathbf{X}}$ and \mathbf{Y} by $\tilde{\mathbf{Y}}$ and results in the parameter vector $\bar{w} = [w; v]$.

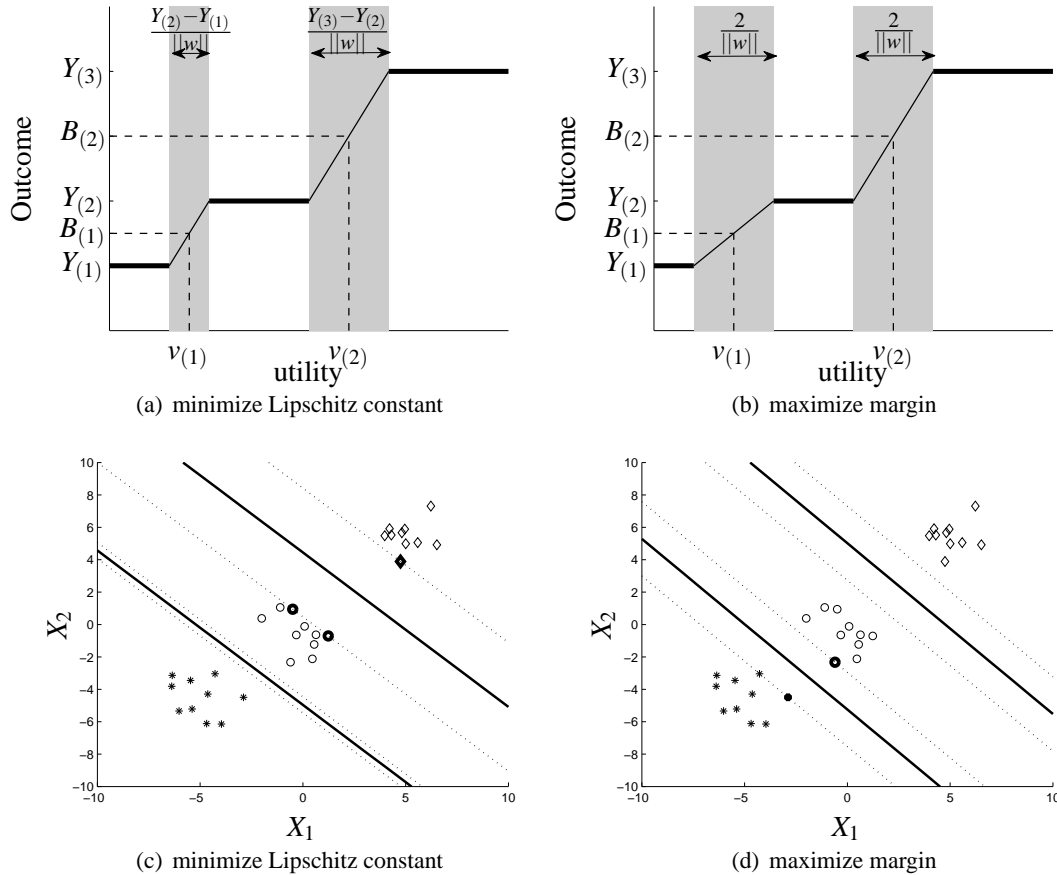


Figure 4: Adaptation of the MINLIP algorithm to ordinal regression: (a) Inclusion of dummy data points with output values B_i intermediate to the observed output values and undefined covariates and utility v_i . All data points are compared with two dummy data points; (b) Comparison with the maximal margin strategy used in standard SVM where the margin is equal between all classes; (c-d) Example with 3 linearly separable cases with outcomes equal to 1 (stars), 2 (circles) and 10 (diamond) respectively. The bold symbols represent the support vectors. In the maximal margin strategy there exists one margin, equal between every two successive classes, which results in a different Lipschitz constant. Using the MINLIP strategy, the Lipschitz smoothness is optimized, resulting in margins which are proportional to the difference in the class labels. Support vectors of the latter method are therefore more likely to be observations of two classes for which the output labels differ the most.

Using the above formulation, the thresholds v as well as the weights w are regularized. Since the motivation for this regularization scheme is not clear, one can formulate the problem explicitly as:

$$\begin{aligned} \min_{w, e, e^*, v} \quad & \|w\|_2 + \gamma \sum_{i=1}^n 1_n^T (e + e^*) \\ \text{s.t.} \quad & \begin{cases} \mathbf{X}w - \mathbf{Q}v + e \geq \mathbf{Y} - \mathbf{Q}\mathbf{B} \\ -\mathbf{X}w + \mathbf{Q}^*v + e^* \geq -\mathbf{Y} + \mathbf{Q}^*\mathbf{B} \\ e \geq 0 \\ e^* \geq 0 \\ \mathbf{M}v \leq 0, \end{cases} \end{aligned} \quad (11)$$

with γ a positive regularization constant, \mathbf{Q} and $\mathbf{Q}^* \in \mathbb{R}^{n \times (k-1)}$ matrices with all elements equal to zero except for positions $\{(i, k_i - 1)\}_{k_i=2}^k$ and $\{(i, k_i)\}_{k_i=1}^{k-1}$ respectively (where k_i represents the index of the output level of observation i), which contain ones. These positions correspond to the dummy data points with which one wishes to compare data points i . Vector $\mathbf{B} \in \mathbb{R}^{k-1}$ contains outcomes corresponding to the thresholds: $\mathbf{B} = [B_{(1)}, \dots, B_{(k-1)}]^T$. Vector v contains all unknown utility function values for the dummy data points $v = [v_{(1)}, \dots, v_{(k-1)}]^T$, and $\mathbf{M} \in \mathbb{R}^{(k-1) \times k}$ gives the first order differences of a vector and is defined as:

$$\mathbf{M} = \begin{bmatrix} 1 & -1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 1 & -1 & 0 & \dots & 0 & 0 \\ \vdots & & & & & & \vdots \\ 0 & \dots & 0 & 0 & 0 & 1 & -1 \end{bmatrix}.$$

The Lagrange dual problem becomes

$$\begin{aligned} \min_{\alpha, \beta} \quad & \frac{1}{2} \alpha^T \mathbf{K} \alpha + \frac{1}{2} \beta^T \mathbf{K} \beta - \alpha^T \mathbf{K} \beta - \alpha^T (\mathbf{Y} - \mathbf{B}^T \mathbf{Q}) + \beta^T (\mathbf{Y} - \mathbf{B}^T \mathbf{Q}^*) \\ \text{s.t.} \quad & \begin{cases} 0_n \leq \alpha \leq \gamma 1_n \\ 0_n \leq \beta \leq \gamma 1_n \\ 0_{k-2} \leq v \\ \mathbf{Q}^T \alpha - \mathbf{Q}^{*T} \beta + \mathbf{M}^T v = 0_{k-1}, \end{cases} \end{aligned}$$

where 1_n and 0_n represent column vectors of size n with all elements equal to 1 and 0 respectively. Solving this explicit formulation is computationally less demanding and faster than solving the implicit problem formulation. We refer to Appendix C for a detailed derivation. The estimated \hat{u} can be evaluated at any point $X^* \in \mathbb{R}^d$ as

$$\hat{u}(X^*) = (\hat{\alpha}^T - \hat{\beta}^T) \mathbf{K}_n(X^*),$$

with $\mathbf{K}_n(X^*)$ defined as before.

4.2 Prediction for Ordinal Regression

A clear advantage of the approach which includes unknown thresholds is that the prediction step becomes very simple. As illustrated in Figure 5, the predictions can be easily obtained from the value of the utility function in comparison with the different threshold values.

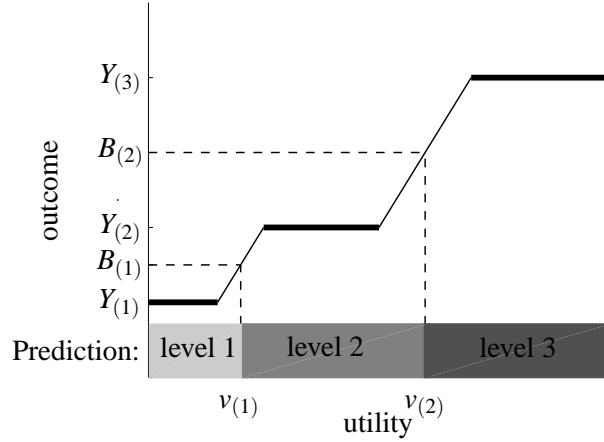


Figure 5: Prediction for ordinal regression. MINLIP for ordinal regression, including unknown thresholds, has the advantage to reduce the prediction step to a simple comparison between the utility of a new observation and the utility of the thresholds. If the utility has a value between threshold $j - 1$ and j , the predicted outcome equals the j^{th} output level.

4.3 Difference with Other Methods

Chu and Keerthi (2005) proposed two SVM based models for ordinal regression. Both methods introduce $k - 1$ thresholds with k the number of ordinal levels in the data. As with SVM classifiers, the margin between two ordinal levels is set to $\frac{2}{\|w\|_2}$. In their first method (EXC) a data point X_i belonging to class Y_i has two slack variables: one relating to the threshold between classes $k_i - 1$ and k_i and a second relating to the threshold between classes k_i and $k_i + 1$. To ensure that the threshold between classes $k_i - 1$ and k_i is smaller than the threshold between classes k_i and $k_i + 1$, $k - 1$ additional constraints are explicitly included. The problem can be written as:

$$\begin{aligned} \min_{w, e, e^*, v} \quad & \|w\|_2 + \gamma \sum_{i=1}^n (e_i + e_i^*) \\ \text{s.t.} \quad & \begin{cases} w^T X_i - v_j + e_i \geq 1 & \forall i = 1, \dots, n; j = \arg \max_j (T_i > v_j) \\ -w^T X_i + v_j + e_i^* \geq 1 & \forall i = 1, \dots, n; j = \arg \min_j (T_i < v_j) \\ e_i \geq 0 & \forall i = 1, \dots, n \\ e_i^* \geq 0 & \forall i = 1, \dots, n \\ v_j \leq v_{j+1} & \forall j = 1, \dots, k - 2. \end{cases} \end{aligned}$$

In their second approach (IMC) the constraints on the thresholds are added implicitly by adding $k - 1$ slack variables, one for each threshold, for every datapoint. The problem can be formulated as follows:

$$\begin{aligned}
 & \min_{w, e, e^*, v} \|w\|_2 + \gamma \sum_{i=1}^n (e_i + e_i^*) \\
 & \text{s.t.} \quad \begin{cases} w^T X_i - v_j + e_i \geq 1 & \forall i = 1, \dots, n; \forall j: T_i > v_j \\ -w^T X_i + v_j + e_i^* \geq 1 & \forall i = 1, \dots, n; \forall j: T_i < v_j \\ e_i \geq 0 & \forall i = 1, \dots, n \\ e_i^* \geq 0 & \forall i = 1, \dots, n. \end{cases}
 \end{aligned}$$

In our method, we adopt the approach of the EXC method concerning slack variables, the method differing in the definition of the margin. Instead of defining an equal margin at every border, the margin between classes k_i and $k_i + 1$ is defined as $\frac{|Y_{(i+1)} - Y_{(i)}|}{\|w\|_2}$.

Remark the similarity between these models and the standard SVMs (Vapnik, 1998) in the binary classification problem (with two classes C_1 and C_2):

$$\begin{aligned}
 & \min_{w, e, e^*, b} \|w\|_2 + \gamma \sum_{i=1}^n (e_i + e_i^*) \\
 & \text{s.t.} \quad \begin{cases} w^T X_i + b + e_i \geq 1 & \forall i \in C_1 \\ -w^T X_i - b + e_i^* \geq 1 & \forall i \in C_2 \\ e_i \geq 0 & \forall i \in C_1 \\ e_i^* \geq 0 & \forall i \in C_2. \end{cases} \quad (12)
 \end{aligned}$$

In case $k = 2$, both EXC and IMC reduce to the model:

$$\begin{aligned}
 & \min_{w, e, e^*, v} \|w\|_2 + \gamma \sum_{i=1}^n (e_i + e_i^*) \\
 & \text{s.t.} \quad \begin{cases} w^T X_i - v + e_i \geq 1 & \forall i \in C_1 \\ -w^T X_i + v + e_i^* \geq 1 & \forall i \in C_2 \\ e_i \geq 0 & \forall i \in C_1 \\ e_i^* \geq 0 & \forall i \in C_2, \end{cases} \quad (13)
 \end{aligned}$$

which equals the model in Equation (12) when the threshold v (note that there is only one threshold in this case) is considered as the constant term. The MINLIP model reduces to:

$$\begin{aligned}
 & \min_{w, e, e^*, v} \|w\|_2 + \gamma \sum_{i=1}^n (e_i + e_i^*) \\
 & \text{s.t.} \quad \begin{cases} w^T X_i - v + e_i \geq Y_i - B & \forall i \in C_1 \\ -w^T X_i + v + e_i^* \geq B - Y_i & \forall i \in C_2 \\ e_i \geq 0 & \forall i \in C_1 \\ e_i^* \geq 0 & \forall i \in C_2, \end{cases} \quad (14)
 \end{aligned}$$

where only one dummy observation (v, B) needs to be introduced. The difference between Equations (12, 13) and Equation (14) lies in the right hand side of the two first inequalities, which is a consequence of the used complexity control. Models (13) and (14) are equivalent up to the choice of the regularization constant.

Chu and Ghahramani (2005) proposed a probabilistic approach to ordinal regression in Gaussian processes (GPOR). They impose a Gaussian process prior distribution on the utility function (called latent function in their work) and employ an appropriate likelihood function for ordinal variables.

Experiments will compare our methods with the Bayesian inference technique of MacKay (1992), using the Laplacian approximation to implement model adaptation. The GPOR approach differs from ours since it uses a Bayesian framework.

5. Transformation Models for Failure Time Data

We now turn our attention to the case where the data originate from a survival study, that is, the dependent variable is essentially a time-to-failure and typically requires specific models and tools to capture its behavior. We will adopt a classical statistical setup, and will show how the techniques as described in Section 3 provide a powerful alternative to the classical statistical (semi-parametric) toolbox.

5.1 Survival Data

The observations are assumed to fit in the following statistical setup, see, for example, Kalbfleisch and Prentice (2002) for a more elaborate introduction. Let $T \in \mathbb{R}^+$ and $X \in \mathbb{R}^d$ be a random variable and random vector respectively, jointly following a probability law characterized by P as classical. The former variable T describes the *time to the event* of interest, and the random vector X taking values in \mathbb{R}^d describes d covariates. Note that in this Section T has the same role as Y in the previous Sections. We assume that no ties will occur in the data in order to keep the explanations as simple as possible. We will consider predictive models where the covariates come in through a linear combination with weights $w \in \mathbb{R}^d$ as before, or $\mathcal{U} = \{u : \mathbb{R}^d \rightarrow \mathbb{R} : u(X) = w^T X, \forall X \in \mathbb{R}^d\}$. A key quantity in survival analysis is the *conditional survival function* $S(t|u(X)) : \mathbb{R}^+ \rightarrow [0, 1]$ defined as

$$S(t|u(X)) = P\left(T > t \mid u(X)\right),$$

denoting the probability of the event occurring past t given the value of the utility function $u(X) = w^T X$. A related quantity to the conditional survival function is the *conditional hazard function* $\lambda : \mathbb{R} \rightarrow \mathbb{R}^+$ defined as

$$\begin{aligned} \lambda(t|u(X)) &= \lim_{\Delta t \rightarrow 0} \frac{P\left(t \leq T < t + \Delta t \mid u(X), T \geq t\right)}{\Delta t} \\ &= \lim_{\Delta t \rightarrow 0} \frac{P\left(t \leq T < t + \Delta t \mid u(X)\right)}{S\left(t \mid u(X)\right)}. \end{aligned}$$

If the derivative $s : \mathbb{R}^+ \rightarrow \mathbb{R}$ with $s(t|u(X)) = \frac{\partial S(t|u(X))}{\partial t}$ exists, one can write $\lambda(t|u(X)) = \frac{s(t|u(X))}{S(t|u(X))}$. The conditional hazard function reflects the instantaneous probability that the event will occur given that the subject already survived beyond time t . Finally, one can make the relation between the hazard λ and the survival function S even more explicit by introducing the *conditional cumulative hazard function* $\Lambda(t|u(X)) = \int_0^t \lambda(r|u(X)) dr$ for $t \geq 0$ such that

$$\Lambda(t|u(X)) = -\ln\left(S(t \mid u(X))\right).$$

The following Subsection enumerates some commonly used (semi-)parametric methods for modelling the survival and hazard functions.

5.2 Transformation Models for Survival Analysis

The Transformation model (see Definition 1) encompasses a broad class of models, including the following classical survival models.

1. Cox' *proportional hazard model* is recovered when one defines $g = h^{-1}$ (if it exists) as $g(z) = \ln(-\ln(z))$. Under the Cox model, the value of the survival function at $t = T$ is

$$S(T, X) = [S_0(T)]^{\exp(-\beta^T X)},$$

where $S_0(t) = \exp(-\Lambda_0(t))$ is called the baseline survival function. Taking $\ln(-\ln(\cdot))$ of both sides in (1) leads to

$$\begin{aligned} \ln(-\ln(S(T, X))) &= \ln(-\ln(S_0(T))) - \beta^T X \\ &= \ln(\Lambda_0(T)) - \beta^T X \\ \Rightarrow \varepsilon &= g(T) - u(X) \\ \Rightarrow T &= h(u(X) + \varepsilon). \end{aligned}$$

Remark that the last transition is only possible if $g(t)$ is invertible. The 'noise terms' are i.i.d. observations from the extreme value distribution $F_\varepsilon(z) = 1 - \exp(-\exp(z))$.

2. The *proportional odds model* is defined as

$$\ln\left(\frac{F(t|X)}{1 - F(t|X)}\right) = \alpha(t) + \beta^T X, \quad (15)$$

with $F(t|X)$ the conditional cumulative distribution function and $\alpha(t)$ a monotonically increasing function. In general the survival function equals $S(t) = 1 - F(t)$, leading together with Equation (15) to

$$\begin{aligned} \ln\left(\frac{1 - S(T|X)}{S(T|X)}\right) &= \alpha(T) + \beta^T X \\ \Rightarrow \varepsilon &= \alpha(T) + u(X) \\ \Rightarrow T &= h(-u(X) + \varepsilon). \end{aligned}$$

Remark that the last transition is only possible if $\alpha(T)$ is invertible.

3. The *accelerated failure time (AFT)* is given when $h(z) = \ln(z)$.

For an extended discussion on the use of the class of transformation models and specific parameterizations of the functions h or g , see, for example, Dabrowska and Doksum (1988), Koenker and Geling (2001), Cheng, Wei, and Ying (1997) and citations.

5.3 Censoring

A typical property of failure time data is the occurrence of censoring. A failure time is called censored when the exact time of failure is not observed. Despite this, censored times do provide relevant information. Define $\mathcal{T}_i = (T_i, \delta_i)$ with δ_i the censoring indicator, capturing all censoring

information: $\delta = 0$ indicates the occurrence of an event at a known failure time (uncensored data point); right, left and interval censoring are indicated by $\delta = 1$, $\delta = 2$ and $\delta = 3$ respectively. Without censoring all possible pairs of datapoints $\{(T_i, T_j)\}_{i \neq j}$ can be used for comparison in Equation (5). The presence of censoring leads to a lack of comparability between certain data points. Let $\Delta(T_i, T_j)$ be a comparability indicator, indicating whether the datapoints i and j are comparable:

$$\Delta(T_i, T_j) = \begin{cases} 0 & \text{if } T_i \text{ and } T_j \text{ are not comparable} \\ 1 & \text{if } T_i \text{ and } T_j \text{ are comparable.} \end{cases}$$

This indicator is defined depending on the censoring types present in the data:

Right censoring occurs when the event of interest did not occur until the last follow-up time. This type of censoring typically occurs at the end of the study period. Although the exact failure time is not known in this case, the failure time is known to be later than the date of last follow-up. In case of right censoring the comparability indicator Δ takes the value 1 for two observations i and j when the observation with the earliest failure time is observed, and zero otherwise:

$$\Delta(T_i, T_j) = \begin{cases} 1 & \text{if } (T_i < T_j \text{ and } \delta_i = 0) \text{ or } (T_j < T_i \text{ and } \delta_j = 0) \\ 0 & \text{otherwise.} \end{cases}$$

Left censoring deals with the case when the failure is known to have happened before a certain time. An example of left censoring arises in case a variable can only be measured when its value is above a certain level. For left censoring, two observations i and j are comparable when the observation with the highest failure time is non-censored and zero otherwise:

$$\Delta(T_i, T_j) = \begin{cases} 1 & \text{if } (T_i < T_j \text{ and } \delta_j = 0) \text{ or } (T_j < T_i \text{ and } \delta_i = 0) \\ 0 & \text{otherwise.} \end{cases}$$

Interval censoring is a combination of the previous two censoring types. In this case the failure time is not known exactly, instead an interval including the failure time is indicated. This type of censoring is often found in medical studies where the patients are subject to regular check up times (Finkelstein, 1986). Whether two observations are comparable or not in case of interval censoring depends on the censoring times \underline{T}_i and \bar{T}_i defining the failure interval for each observation i : $T_i \in [\underline{T}_i, \bar{T}_i]$. For uncensored observations, the failure interval reduces to one time, namely the failure time $T_i = \underline{T}_i = \bar{T}_i$. The comparability indicator is defined as:

$$\Delta(T_i, T_j) = \begin{cases} 1 & \text{if } \bar{T}_i < \underline{T}_j \text{ or } \bar{T}_j < \underline{T}_i \\ 0 & \text{otherwise.} \end{cases}$$

In case the data consists of data points with different types of censoring, the comparability indicator is defined as follows. In the most general case, the failure time T_i is considered to be an element of the interval $[\underline{T}_i, \bar{T}_i]$. For right censored data points the right edge of the interval equals infinity, whereas for left censored observation the left edge of the interval equals zero. The comparability indicator is then defined as:

$$\Delta(\mathcal{T}_i, \mathcal{T}_j) = \begin{cases} 1 & \text{if } \bar{T}_i < \underline{T}_j \text{ or } \bar{T}_j < \underline{T}_i \\ 0 & \text{otherwise.} \end{cases}$$

More information on censoring can be found in Andersen et al. (1993), Elandt-Johnson and Johnson (1980), Harrell (2001), Kalbfleisch and Prentice (2002) and Miller (1981).

Standard statistical methods for modelling survival data obtain parameter estimates by maximizing a (partial) likelihood with regard to these parameters. This likelihood depends on the ranking of the failure times. In the presence of right censoring, this ranking can uniquely be defined and estimates for the parameters can be obtained. However, in the presence of interval censoring, a unique ranking of the failure of all instances is not always possible. Peto (1972) and Satten (1996) among others, suggested extensions of the proportional hazard model where censoring is not restricted to right censoring. However, estimation of the parameters in these cases remain difficult. In the next section we will illustrate that MINLIP can be easily adapted for right, left, interval censoring and combined censoring schemes. However, we first need an appropriate measure of concordance equivalent to Equation (3). Therefore, we resort to the concordance index as described by Harrell et al. (1984) and Harrell (2001).

Definition 4 (Concordance Index) *The concordance index (c-index) is a measure of association between the predicted and observed failures in case of censored data. The c-index equals the ratio of concordant to comparable pairs of data points. Two observations i and j are comparable if their relative order in survival time is known. A pair of observations i and j is concordant if they are comparable and the observation with the lowest failure time also has the lowest score for the utility function $u(X)$. Formally, the observation based c-index of a model generating predictions $u(X_i)$ for data X_i from a data set $\mathcal{D} = \{(X_i, Y_i, \delta_i)\}_{i=1}^n$ can be expressed as*

$$C_n(u) = \frac{\sum_{i \neq j} \Delta(\mathcal{T}_i, \mathcal{T}_j) I[(u(X_j) - u(X_i))(Y_j - Y_i) > 0]}{\sum_{i \neq j} \Delta(\mathcal{T}_i, \mathcal{T}_j)}.$$

This index is an estimate probability of concordance between predicted and observed survival, with c-index = 0.5 for random predictions and c-index = 1 for a perfectly concordant model. Without censoring, this definition is exactly equal to the concordance as defined in Equation (2).

5.4 Modifications to MINLIP

This section describes how the standard MINLIP model can be extended towards failure time data including the handling of censored data. Therefore, Equation (8) is adapted to include censored data. In particular, the matrix \mathbf{D} needs to be changed in order to allow for pairs of data points not to be comparable. Let $\mathbf{R} \in \mathbb{R}^{(n-1) \times (n-1)}$ be defined as the diagonal matrix with $\mathbf{R}_{ii} = \Delta(Z_i, Z_{i+1}), \forall i = 1, \dots, n-1$. The matrix \mathbf{D} , representing the datapoints to be compared, is adapted for censoring according to:

$$\mathbf{D}_c = \mathbf{R}\mathbf{D},$$

where Δ is defined as in Section 5.3, resulting in multiple rows with only zero entries in the matrix \mathbf{D}_c . For computational convenience these rows can be left out. It is seen that issues concerning the

type(s) of censoring in the data are easily dealt with by using the comparability indicator. In the remainder of this paper we will restrict our attention to right censored data.

The learning objective can now be formalized as

$$\begin{aligned} \min_{w,e} \quad & \frac{1}{2}w^T w + \gamma \|e\|_1 \\ \text{s.t.} \quad & \mathbf{D}_c(\Phi w + e) \geq \mathbf{D}_c \mathbf{T}, \end{aligned} \tag{16}$$

where $\|e\|_1 = \sum_{i=1}^n |e_i|$ and $\mathbf{T} = [T_{(1)}, T_{(2)}, \dots, T_{(n)}]^T$ is a vector containing all failure times, censored or not. As in Section 3, the dual of this optimization problem becomes

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2}\alpha^T \mathbf{D}_c \mathbf{K} \mathbf{D}_c^T \alpha - \alpha^T \mathbf{D}_c \mathbf{T} \\ \text{s.t.} \quad & \begin{cases} -\gamma \mathbf{1}_n \leq \mathbf{D}_c^T \alpha \leq \gamma \mathbf{1}_n \\ \alpha \geq 0_{n-1}, \end{cases} \end{aligned}$$

Given the solution $\hat{\alpha}$, the predicted utility can be calculated for a new point X^* as

$$u(X^*) = \hat{\alpha}^T \mathbf{D}_c \mathbf{K}_n(X^*),$$

with $\mathbf{K}_n(X^*) = [K(X^*, X_1) \dots K(X^*, X_n)]^T \in \mathbb{R}^n$. Since the censoring mechanism can be handled by a proper choice of \mathbf{D}_c , it is not too difficult to extend the formulations of Subsection 3.5 as well.

5.5 Prediction with Transformation Models

The prediction step in survival analysis, refers to the estimation of survival and hazard functions rather than the estimation of the failure time itself. The proportional hazard model estimates these functions, by assuming that a baseline hazard function exists; the covariates changing the hazard only proportionally. The baseline hazard function is estimated using the Breslow estimator of the cumulative baseline hazard (Breslow, 1974).

In our setting, the cumulative distribution function (cdf), can be estimated, after estimation of the utility, as follows. The time axis is divided in k equidistant time intervals $[t_{l-1}, t_l], \forall l = 2, \dots, k$. For each observation in the set $\{u_i, T_i, \delta_i\}_{i=1}^n$, the outcome in each time interval is defined as:

$$Y_{il} = \begin{cases} 0 & \text{if } T_i > t_l \\ 1 & \text{if } T_i \leq t_l \text{ and } \delta_i = 0. \end{cases}$$

Remark that censored observations are not considered at times later than the censoring time. Using a monotone least squares support vector regression model (Pelckmans et al., 2005a) with a Gaussian kernel, or another monotonic regression model, the utility and the time interval number l as inputs and Y_{il} as output, the cdf $\hat{F}(u_i, l)$ is estimated. The survival function is found as $\hat{S}(u_i, l) = 1 - \hat{F}(u_i, l)$. The hazard function is then found as

$$\hat{\lambda}(u_i, l) = \frac{\hat{F}(u_i, l+1) - \hat{F}(u_i, l)}{t_{l+1} - t_l}, \forall l = 1, \dots, k-1.$$

Remark the analogy with the partial logistic artificial neural network approach to the survival problem proposed by Biganzoli et al. (1998). However, since the latter can not be seen as a transformation model, data replication is necessary even when one is only interested in risk groups. Thanks to the two-step approach of transformation models, data replication can be avoided.

6. Application Studies

This final Section describes some experiments to illustrate the use of the presented method. In a first Subsection, 3 artificial examples will illustrate how transformation models are used within the ranking, ordinal regression and survival setting (see also Table 1). A first real life application illustrates the use of MINLIP for ordinal data. We use 6 benchmark data sets and compare the performance of MINLIP with EXC and IMC as proposed in Chu and Keerthi (2005) and GPOR as proposed in Chu and Ghahramani (2005). The last two examples concern survival data, one with micro-array data (data also used in Bøvelstad et al. 2007) and one with clinical data (Schumacher et al., 1994).

Unless stated otherwise, 10-fold cross-validation was used for model selection. For every kernel and regularization parameter to be tuned a grid of values was searched and the combination of parameter values yielding the lowest cross-validation error or highest cross-validation performance was selected. In the first example the mean absolute error between predicted and observed output levels was used as model selection criterion since prediction is relevant in this case. For both survival examples, the cross-validation concordance index was used as model selection criterion since the main interest lies in the ranking of the patients.

6.1 Artificial Examples

This section illustrates the different steps needed to obtain the desired output for ranking, regression and survival problems, using artificial examples. Together with Table 1, this Subsection illustrates the different tasks considered in the paper.

6.1.1 RANKING

In this first example, we consider the ranks of 150 cyclists in 9 different races. Using the ranks of 100 out of these cyclists in a 10th race, we want to predict the rank of the remaining 50. Additional information includes: age, weight and condition score. The outcome is defined as the ranking given by a weighted sum of the ranking in the 9 races, age, weight and condition score. Weights are drawn from a uniform distribution on the unit interval. The ranking on the previous races are numbers from 1 to 100, all other variables are drawn from a standard normal distribution.

A first step in all transformation models is to estimate the utility function u . Using Equation (8) with a linear kernel and 5-fold cross validation with the concordance index (ranking criterion) as a model selection criterion, a concordance index of 0.98 on the test set was obtained. The predicted ranking corresponds very well with the observed ranking (see Figure 6). Since one is only interested in ranking the cyclists, the value of the utility is irrelevant. Additionally, one is not interested in estimating the transformation function h .

6.1.2 ORDINAL REGRESSION

In a second artificial example, consider the scenario in which one wishes to divide students into 3 groups: bad, average and good student. For this task, the grades on 10 courses for 150 students are available. The outcome depends on the average grade. A bad, average and good student has an average grade below 55%, between 55% and 65% and above 65% respectively. The results on 100 students will be used for training, the remaining students will be used for testing.

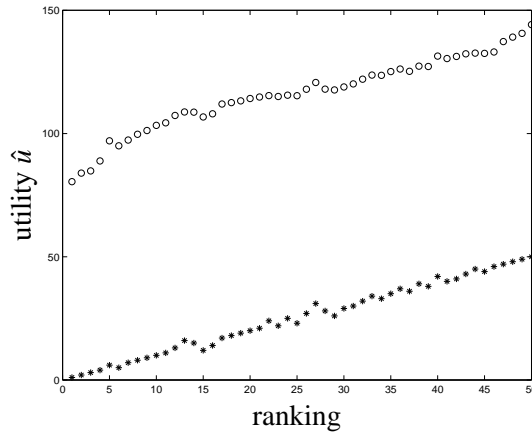


Figure 6: Artificial example illustrating the use of the MINLIP transformation model for the ranking setting. The estimated utility of the test observations are denoted by the circles. The value of the utility does not correspond to the ranking. However, the rank of the estimated utility (denoted by stars) are a good prediction of the observed rank.

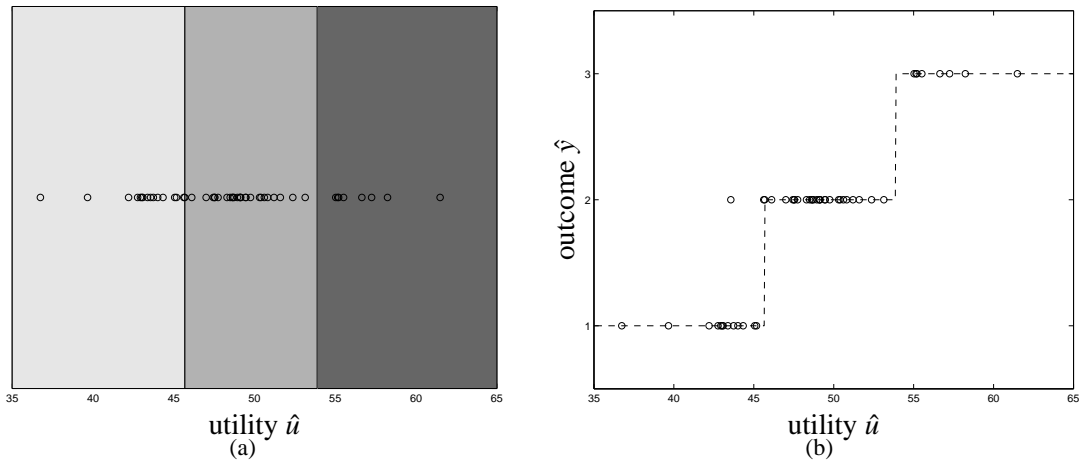


Figure 7: Artificial example illustrating the use of the MINLIP transformation model for ordinal regression. (a) The estimated utility of the test observations are denoted by the circles. The MINLIP model for ordinal regression results in an estimate of the utility function and threshold values. Students with the lowest utility (less than the first threshold) are predicted to be bad students (light grey). Students with a utility between both thresholds (medium grey) are estimated to be average students and students with a utility higher than the second threshold (dark grey) are predicted to be good students. (b) Illustration of the transformation function h (dashed line).

In a first step, all students are ranked according to the available data, namely their grades on 10 courses. Using Equation (11) with the concordance index as model selection criterion and a linear kernel, a concordance of 0.99 between the utility and the outcome on the test set is obtained. In addition to an estimate of the utility, the MINLIP model for ordinal regression gives threshold values which can be used to predict the outcome of new observations (see Figure 7). Since, the MINLIP model generates these thresholds, no additional model is needed to obtain the transformation function.

6.1.3 SURVIVAL ANALYSIS

In a last artificial example, a randomized trial is simulated. Assume 150 patients are randomly divided into 2 treatment groups. Additionally, the age (drawn from a standard normal distribution) of the patients is known. The survival time is known for 100 patients. For the first treatment arm, the survival time has a Weibull distribution with parameters 1 and 0.5. For the second treatment arm, the survival time is Weibull distributed with parameters 4 and 5. Using the information on age, treatment arm and survival on 100 patients, one would like to predict the survival for the remaining 50 patients. Assuming that the age is irrelevant for the survival, the treatment will be the only important factor in predicting the patients' survival.

As with the two previous examples, the MINLIP model is used to estimate the utility of the patients. Using a linear kernel and 5-fold cross validation in comparison with the concordance index as model selection criterion, a c-index of 0.70 is obtained on the test set. Figure 8 illustrates that the MINLIP model is able to divide the group of test patients into two groups with a significant different survival ($p=0.03$, logrank test). The first part of the transformation model obtained a nice result. However, in survival analysis, additional information can be provided when performing the second part of the transformation model, namely estimating the transformation function. Applying the method as explained in Section 5.5, the estimated survival curves for all patients are calculated (Figure 9). One clearly notices two distinct survival groups. The grey and black survival curves correspond to patients in the first and second treatment arm, respectively. The true survival function for the first and second treatment are illustrated in thick black and grey lines, respectively.

6.2 Ordinal Regression

At first 6 regression data sets³ were converted to ordinal regression data sets as follows. The data sets were divided into 20 folds with 10 equal-frequency bins. The output value for each bin was set to the average output within the bin. The performance of the MINLIP model was compared with two methods described in Chu and Keerthi (2005) (see Table 2). Both of these methods optimize multiple thresholds to define parallel discriminant hyperplanes for the ordinal levels. The first method (EXC) explicitly imposes the ordering of the thresholds, whereas this is done only implicitly in the second method (IMC). Tuning of the Gaussian kernel parameter and the regularization parameter was performed with 10-fold cross-validation on an exponential grid using mean absolute error as model selection criterion. After an initial search, a finer search was performed in the neighborhood of the initial optimum. Results of the GPOR method are reported for comparison. The GPOR has a lower mean zero-one error on small data sets. For larger data sets and for mean absolute errors, it performs less. The IMC method has the disadvantage that large QP problems need to be solved for

3. Data are available at <http://www.liacc.up.pu/~ltorgo/Regression/DataSets.html>.

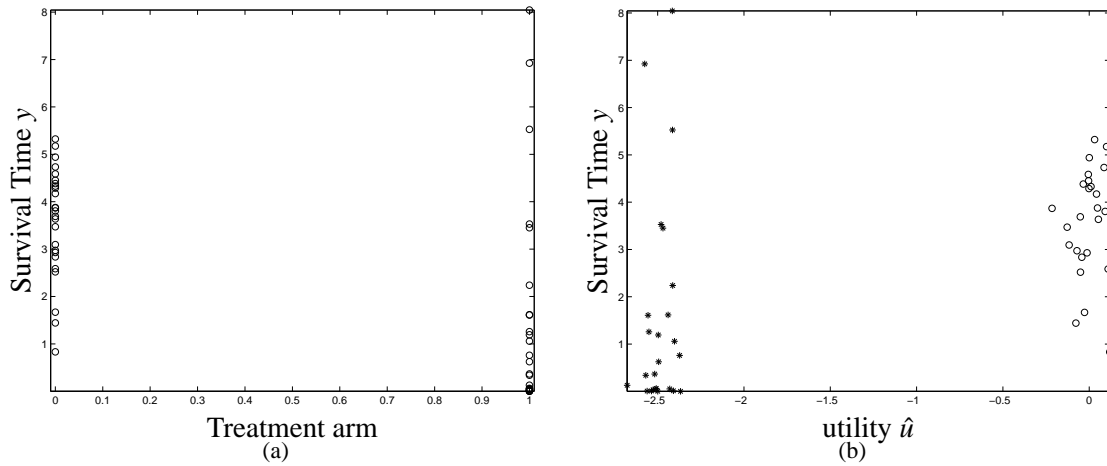


Figure 8: Artificial example illustrating the use of the MINLIP transformation model for survival analysis. (a) Survival time as a function of the treatment arm. Patients receiving the first treatment survive longer in general. The second treatment results in lower survival times. However, some patients have extremely large survival times. (b) Survival time versus estimated utility. The circles and the stars denote the first and second treatment arm, respectively. The utility is able to group the patients according to the relevant variable treatment (see the clear separation between circles and stars).

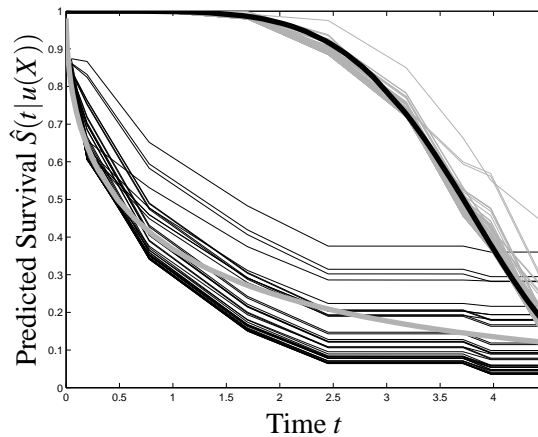


Figure 9: Artificial example illustrating the use of the MINLIP transformation model for survival analysis: illustration of the reconstruction step. For each patient, the survival curve is calculated using the method discussed in Section 5.5. Grey and black curves represent patients from the first and second treatment arm, respectively. The true survival curve for the first and second treatment, are illustrated in thick black and grey lines. One clearly notices two distinct survival groups, corresponding to the treatment groups.

data set	mean zero-one error			
	minlip	exc	imc	gpor
pyrimidines	0.74±0.07	0.79±0.08	0.75±0.09	0.73±0.09
triazines	0.86±0.05	0.86±0.04	0.87±0.04	0.86±0.03
Wisconsin	0.89±0.02	0.89±0.03	0.89±0.02	0.85±0.03 ***
machine CPU	0.65±0.04	0.65±0.06	0.64±0.04	0.80±0.09***
auto MPG	0.58±0.03	0.58±0.03	/	0.68±0.08***
Boston housing	0.57±0.04	0.57±0.04	/	0.61±0.03***
data set	mean absolute error			
	minlip	exc	imc	gpor
pyrimidines	0.05±0.01	0.06±0.01	0.05±0.01	0.06±0.01
triazines	0.12±0.06	0.10±0.01 *	0.10±0.00 ***	0.16±0.02***
Wisconsin	29.56±2.18	28.80±1.17	28.23±1.28	34.52±4.27***
machine CPU	29.41±4.27	30.57±6.94	30.48±4.10	157.00±124.67***
auto MPG	2.08±0.16	2.05±0.18	/	4.13±1.13***
Boston housing	2.49±0.27	2.33±0.25 *	/	2.97±0.29***

Table 2: Test results of MINLIP, EXC and IMC using a Gaussian kernel and GPOR. The targets of the data sets were discretized by 10 equal-frequency bins. The output value for each bin is set to the average output within the bin. The results are averaged over 20 trials. The best performing model is indicated in bold. Significant differences as calculated by Wilcoxon’s signed rank sum test between the EXC, IMC, GPOR and the MINLIP (reference) model are indicated with * ($p < 0.05$), ** ($p < 0.01$) or *** ($p < 0.001$).

growing training samples, requiring more computational time. The MINLIP method makes a nice trade-off between computational time and performance.

6.3 Failure Time Data: Micro-array Studies

The MINLIP technique is derived from machine learning techniques as SVMs, techniques which are shown to be especially useful to handle high-dimensional data sets. We therefore test the performance of MINLIP on 3 micro-array data sets.

In this example we compare the performance of model (16) (MINLIP) and linear extension as in Equation (10) (MINLIP_{L1}) with 5 other methods which are discussed and implemented by Bøvelstad et al. (2007): principal components regression (PCR), supervised principal components regression (SPCR), partial least squares regression (PLS) and two penalized Cox regression models (Cox, 1972): ridge regression (Cox_{L2}) and L1 regularization (Cox_{L1}). The PCR method uses principal component analysis to select n_λ principal components which account for as much variation in the gene expression profiles as possible. The selected principal components are then used as covariates in a Cox regression model (see Hastie, Tibshirani, and Friedman, 2001, Chapter 3.4.4). The SPCR (Bair and Tibshirani, 2004; Bair, Hastie, Debnath, and Tibshirani, 2006) method first selects a subset of genes which are correlated with survival by using univariate selection and then applies PCR to this subset. The standard PLS method performs regression of the outcome using n_λ components which are a linear combination of the original covariates (Martens and Næs, 1989). The application of PLS

to the Cox model is not straightforward since the PLS algorithm assumes a linear relation between outcome and covariates. See Nygård et al. (2008) for a detailed description of the method.

Three different micro-array data sets are used in this experiment:

The Dutch Breast Cancer Data (DBCD) from van Houwelingen et al. (2006) is a subset of the data from van de Vijver et al. (2002) and contains information on 4919 gene expression levels of a consecutive series of 295 women with breast cancer from the fresh-frozen-tissue bank of the Netherlands Cancer Institute. All 295 tumors were primary invasive breast carcinoma less than 5 cm in diameter. The women were 52 years or younger. The diagnosis was made between 1984 and 1995 and there was no previous history of cancer, except non-melanoma skin cancer. In 79 (26.78%) patients distant metastases were noted within the study period. The median follow-up was 6.7 years (range 0.05-18.3).

The DLBCL data from Rosenwald et al. (2002) contains data on 240 patients with diffuse large-B-cell lymphoma. The data consist of 7399 gene expression measurements. The median follow-up time was 2.8 years and 58% of the patients died during the study period.

The Norway/Stanford breast cancer data (NSBCD) from Sørli et al. (2003) contains gene expression measurements from 115 women with breast cancer. 549 intrinsic genes introduced in Sørli et al. (2003) were used. Missing values were previously imputed using 10-nearest neighbor imputation (Bøvelstad et al., 2007). 38 (33%) patients experienced an event.

Figure 10 summarizes performances C_n^u on all methods for 100 different randomizations between training and test sets (2/3 training, 1/3 test). In the right panels of Figure 10 the time dependent receiver operator characteristics (TDROC) (Heagerty, Lumley, and Pepe, 2000) are shown. The left panel illustrates the concordance index. The performance of the MINLIP model is better or comparable to the best of the other tested models.

6.4 Failure Time Data: Cancer Study

In this last example, we investigate the ability of the MINLIP model to estimate how the different covariates influence the survival time. We use the German Breast Cancer Study Group data⁴ (Schumacher et al., 1994), containing information on 686 patients and 8 variables. Available variables are: hormonal treatment, age, menopausal status, tumor size, tumor grade, the number of positive lymph nodes, the progesterone receptor (fmol) and the estrogen receptor (fmol). 299 (43.6%) patients had a breast cancer related event within the study period, leaving all other patients with a right censored failure time. The data set was randomly divided in training and test set (2/3 versus 1/3).

Since medical data are typically not highly non-linear, we use a componentwise polynomial kernel

$$K(X, Z) = \sum_{p=1}^d (\tau + X^{pT} Z^p)^2, \tau \geq 0,$$

with d the number of variables and X^p the p^{th} covariate, to model non-linearities. Model selection is done by 10-fold cross-validation with the concordance index as model selection criterion.

4. Data can be found at <http://www.blackwellpublishers.com/rss/Volumes/A162p1.htm>.

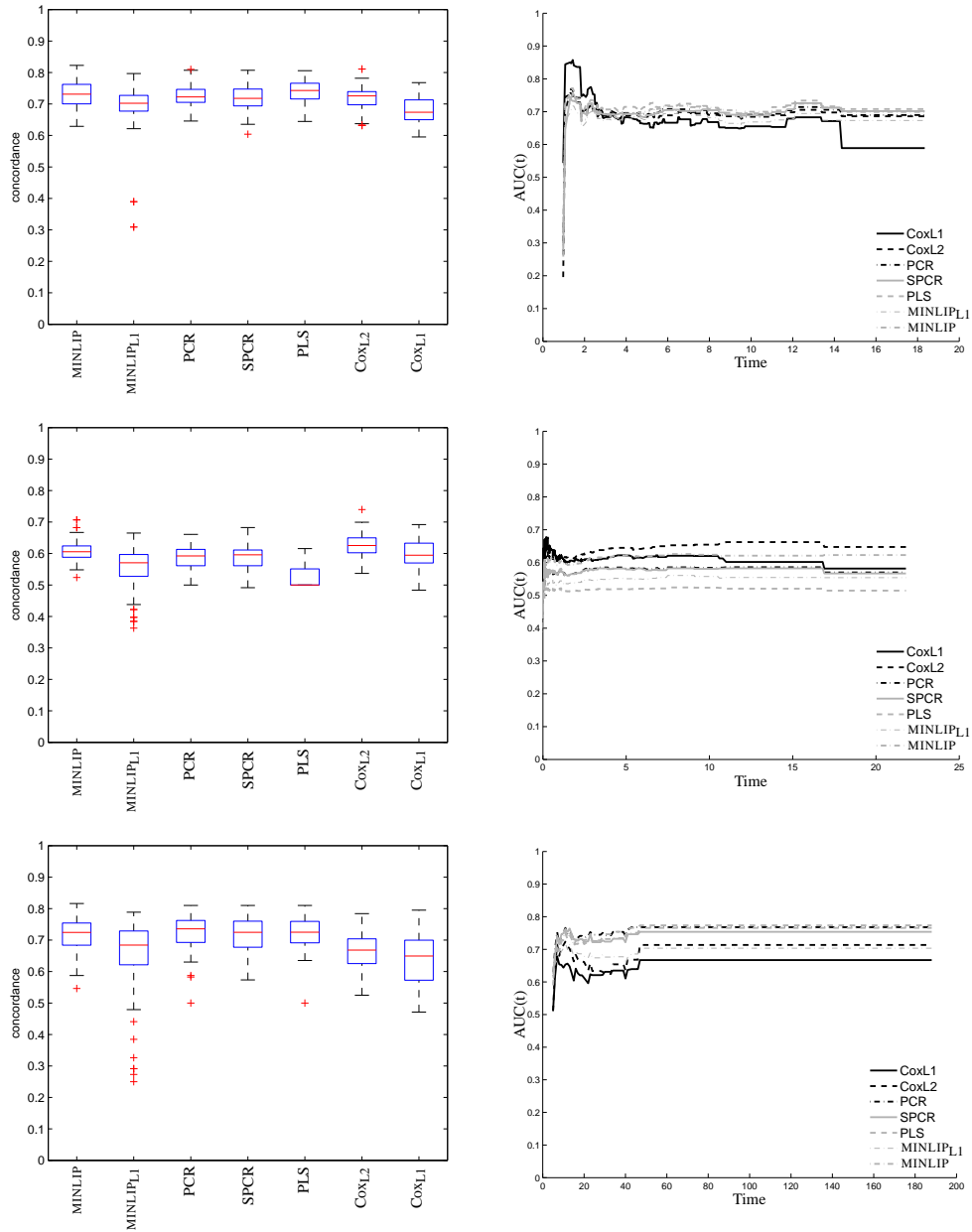


Figure 10: Concordance (left) and time dependent receiver operating characteristic curve (TDROC) (right) on the test set for three micro-array survival data sets (top: DBCD, middle: DL-BCL, bottom: NSBCD). The MINLIP model obtains a performance which is slightly higher or comparable to the other tested models.

We compare our results with Cox' proportional hazard model. However, the standard Cox model (Cox, 1972) assumes linearity in the covariates, implying that for a continuous variable as age for example, the risk ratio between two patients aged 45 and 50 is the same as the risk ratio between two patients aged 75 and 80. To allow for non-linearities in the effects of the covariates on the hazard, the functional forms of the covariates were estimated using penalized smoothing splines (Eilers and Marx, 1996; Hurvich, Simonoff, and Tsai, 1998). In this method, a comparative small set of basis functions is fit and a likelihood penalizing the integrated second derivatives (see Therneau and Grambsch, 2000, Section 5.5) is used to estimate the coefficients. Akaike's information criterion ($AIC = \log \text{likelihood} - \text{degrees of freedom}$) is used to select the degrees of freedom for each term.

Figures 11 and 12 show the estimated covariate effects for Cox regression with penalized splines and MINLIP, respectively. Remark that in Figure 11 the estimates are inversely related with the survival time, whereas in Figure 12 the estimates are related with the survival time itself. Cox' model predicts a decreasing risk for relapse for older patients, up to the age of 40, whereafter the risk increases slowly; for tumors up to 20mm the risk for relapse increases with size, with a threshold effect for larger tumors; the number of positive lymph nodes is inversely related with survival and larger values for the progesterone and estrogen receptors are related with longer survival. All conclusions of the covariate effects agree with what is known from literature (Fisher et al., 1983; Lamy et al., 2002; Pichon et al., 1980; Verschraegen et al., 2005). The MINLIP model estimates a higher survival time for older patients, up to the age of 65, whereafter the survival time drops again. According to this model, a larger tumor, a higher number of positive lymph nodes and a lower progesterone and estrogen receptor level result in lower survival times and thus a higher risk for relapse. Cox' model with penalized smoothing splines obtains a concordance on the test set equal to 0.6715, while the MINLIP model obtains a performance of 0.6857.

Figure 13 illustrates the ability of the models to generate prognostic indices. In clinical practice one is interested in groups of patients with low/high risk for the event to occur. Therefore the median value of the model output is used as a threshold to divide the test set into two groups: one group including patients with an estimated risk lower than the average and a second group with an estimated risk higher than the average. Kaplan-Meier curves and 95%-confidence intervals are plotted in Figure 13. The logrank test χ^2 value is 20.4281 and 29.6984 for Cox and MINLIP respectively. The latter method results in a better split between low and high risk patients.

7. Conclusions

This paper studied a machine learning approach for finding transformation models. Such models are found useful in a context of ordinal regression and survival analysis, and relate directly to commonly used risk measures as the area under the curve and others. The derivations go along the same lines as used for support vector machines, except for replacing the notion of (pairwise) margin with a Lipschitz smoothness condition on the transformation function. The presented learner finds a (non-linear) non-parametric transformation model by solving a convex Quadratic Program. Extensions towards tasks where transformation models provide only a (good) approximation (agnostic case), ordinal regression and survival analysis are given. Experiments on ordinal regression and survival analysis, on both clinical and high dimensional data sets, illustrate the use of the proposed method.

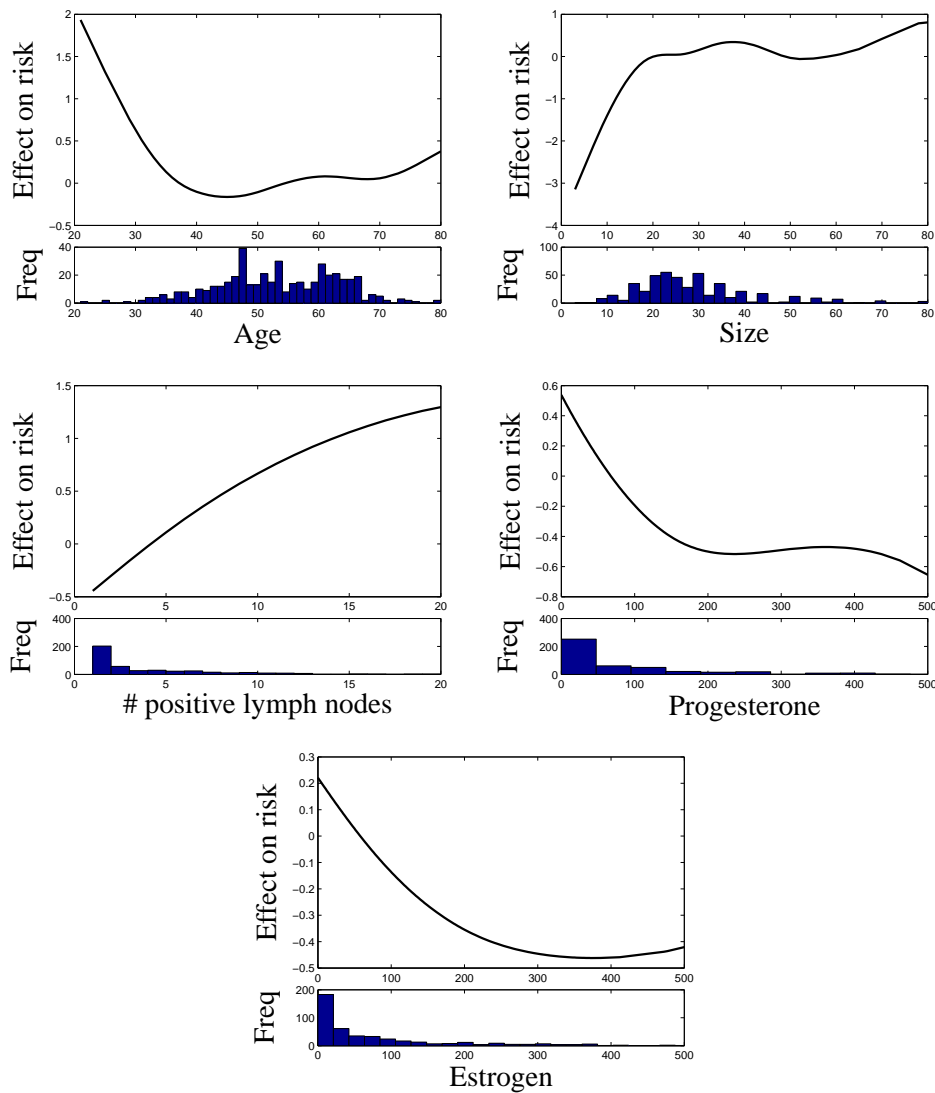


Figure 11: Estimation of the covariate effects on the risk of relapse (remark the difference with Figure 12) with smoothing splines within Cox' proportional hazard model and histograms of the variables. The estimated effects are inversely related with the survival time. The model estimates a lower chance for relapse for older patients up to the age of 40, whereafter the risk increases again, albeit slowly. The chance for relapse increases for larger tumors until a size of 20mm, whereafter the chance remains fairly constant. For common values of the number of positive lymph nodes and receptors, the risk increases for larger/lower values respectively. Conclusions drawn by the model agree with what is known from literature.

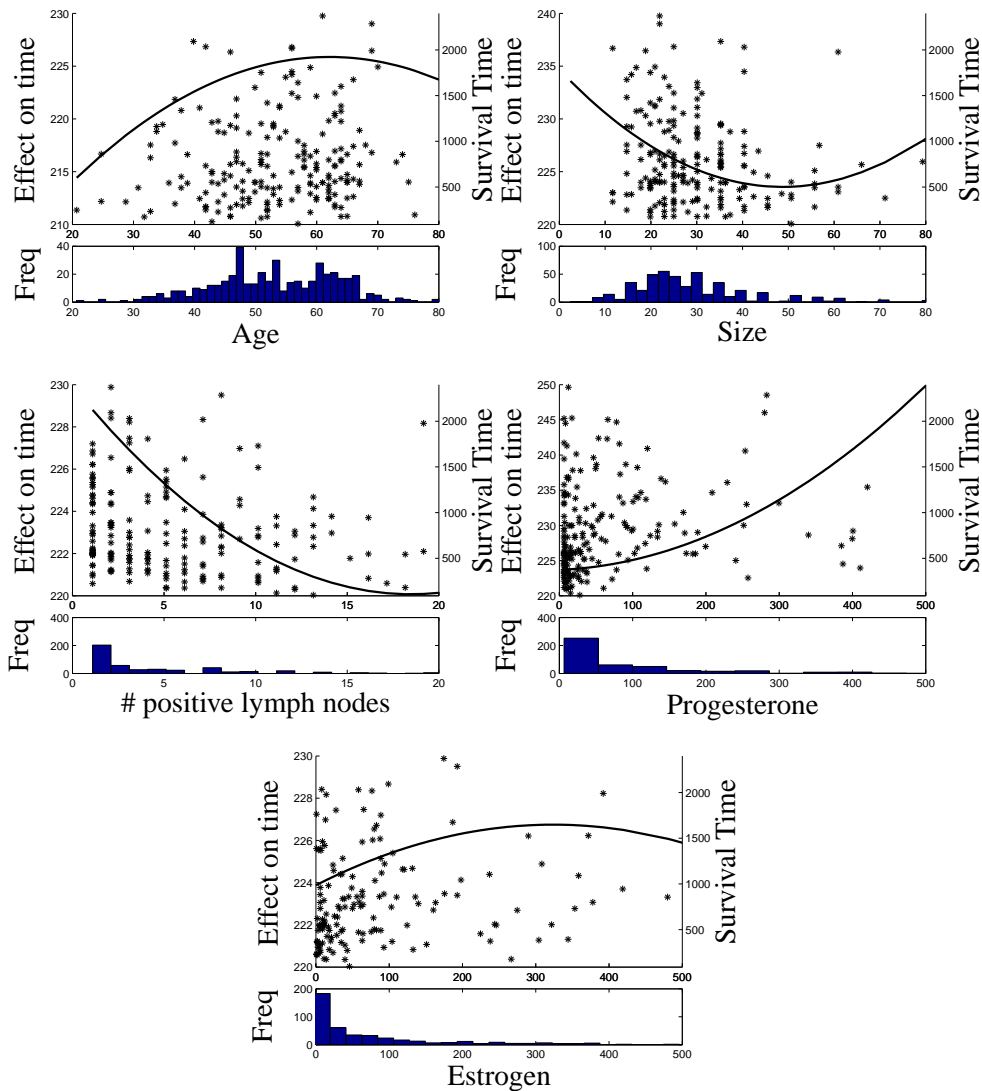


Figure 12: Estimation of covariate effects on survival time (remark the difference with Figure 11) with the MINLIP model ($C_n(u)$ was used for model selection) and histograms of the variables. The stars indicate the observed failure times for breast cancer related events. The estimated covariate effects are directly related with the survival time. The MINLIP model estimates the covariate effects as follows: the estimated survival time increases with age until the age of 65, whereafter the survival time drops slightly. The larger the tumor, the higher the number of positive lymph nodes, the lower the expression of the receptors, the lower the estimated survival time is. Conclusions drawn by the model agree with what is known from literature.

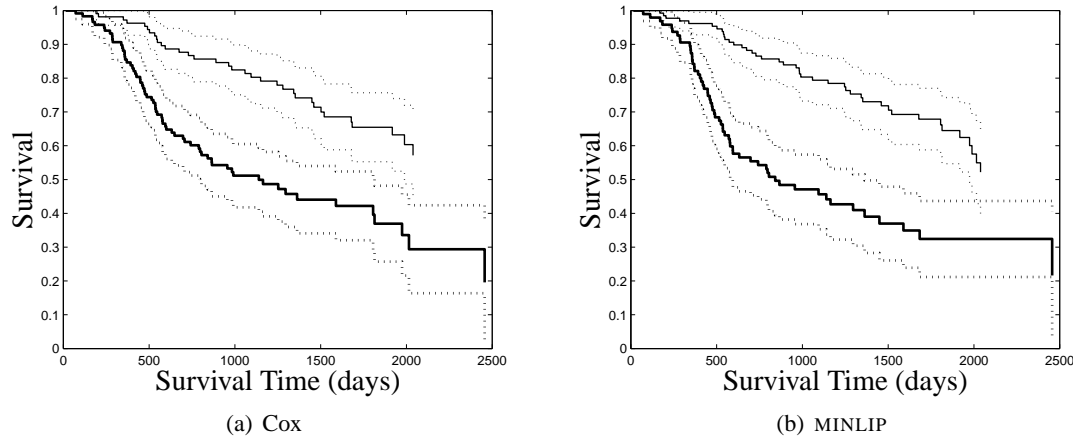


Figure 13: The use of Cox' and MINLIP model as a prognostic index. The output of both models is used to divide the test set into two groups, one with high and one with low risk for relapse. The threshold between both groups is defined as the median value of the model's output. Kaplan-Meier curves and 95% confidence intervals are shown for each group. The spread in the survival curves is broader for the MINLIP model, which is confirmed by a larger value of the log rank test statistic.

Acknowledgments

We thank the editor and the reviewers for their helpful comments and suggestions. This research is supported by Research Council KUL: GOA AMBioRICS, GOA MaNet, CoE EF/05/006, IDO 05/010 IOF-KP06/11, IOF-SCORES4CHEM, several PhD, postdoc end fellow grants; Flemish Government: FWO: PhD and postdoc grants, G.0407.02, G.0360.05, G.0519.06, FWO-G.0321.06, G.0341.07, projects G.0452.04, G.0499.04, G.0211.05, G.0226.06, G.0302.07; IWT: PhD Grants, McKnow-E, Eureka-Flite; Belgian Federal Science Policy Office: IUAP P6/04 (DYSCO, 'Dynamical systems, control and optimization', 2007-2011); EU: FP6-2002-LIFESCI-HEALTH 503094, IST-2004-27214, FP6-MC-RTN-035801; Prodex-8 C90242; EU: ERNSI; the European Research Council under the Seventh Framework Program and within Advanced Grant no. 247035 "Systems and Signals Tools for Estimation and Analysis of Mathematical Models in Endocrinology and Neurology". V. Van Belle is supported by a grant from the IWT. K. Pelckmans is an associated professor/researcher ('forskarassistent') at Uppsala University, Sweden at the department of Information Technology, division of SysCon. S. Van Huffel is a full professor and J.A.K. Suykens is a professor at the Katholieke Universiteit Leuven, Belgium.

Appendix A. Consistency and Identifiability

This first Appendix deals with the issues of consistency and identifiability of the proposed method. We study the question under what conditions MINLIP is consistent, that is, if we have enough data-points at our disposal, would the estimate \hat{w} converge to the desired parameter vector?

Assume that any observation $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ would obey the relation

$$Y = h_0(w_0^T X),$$

where we refer to the (fixed but unknown) vector $w_0 \in \mathbb{R}^d$ as to the 'true' parameters, and to the (fixed but unknown) monotonically increasing function $h_0 : \mathbb{R} \rightarrow \mathbb{R}$ as the 'true' transformation function. We will focus on estimating the vector of parameters w_0 , recovery of h_0 may be done in a second stage. Note that assuming that a finite value L_0 exists, together with an observation (X, Y) where this Lipschitz constant is met, is sufficient for consistency. We will fix $\|w_0\|_2 = 1$ to avoid the obvious identifiability issue, namely that for any strictly positive constant $\alpha > 0$, the system $Y = h_{0,\alpha}(\alpha w_0^T X)$ is not distinguishable from (6) when $h_{0,\alpha}(Z) \triangleq h_0(\frac{Z}{\alpha})$ for any $Z \in \mathbb{R}$.

Let the set of all (possibly an infinite number) observations $\{(X, Y)\} \subset \mathbb{R}^d \times \mathbb{R}$ obeying the system (6) be denoted as \mathcal{D} . We consider that this set is an ε -non-degenerate set, which is defined as follows

Definition 5 (An ε -non-degenerate Set) *Let $\varepsilon > 0$ be any (arbitrarily small) constant. We say that a set $\mathcal{D} = \{(X, Y)\} \subseteq \mathbb{R}^d \times \mathbb{R}$ is ε -non-degenerate if for any observation $(X, Y) \in \mathcal{D}$, and for any vector $v \in \mathbb{R}^d$, one has an observation $(X', Y') \in \mathcal{D}$ different from (X, Y) such that $\|X - X'\|_2 \leq \varepsilon$ so that*

$$v^T (X - X') \geq 0.$$

This requirement can be relaxed as it only has to hold for the point (X, Y) where the Lipschitz condition is met. In addition, we assume that h_0 is (L_0, a) -Lipschitz on this set $\mathcal{D} \subseteq \mathbb{R}$.

Definition 6 (h_0 is (L_0, a) -Lipschitz on $\mathcal{D}' \subseteq \mathbb{R}$) *The monotonically increasing function h_0 is said to be (L_0, a) -Lipschitz on $\mathcal{D} \subseteq \mathbb{R}$ if (1) h_0 is Lipschitz smooth with Lipschitz constant L_0 for all pairs Z, Z^* , with $Z \geq Z^*$:*

$$h_0(Z) - h_0(Z^*) \leq L_0(Z - Z^*),$$

and (2) there exists a pair $Z, Z' \in \mathcal{D}$ with $Z > Z'$ where the Lipschitz constant is met:

$$h_0(Z) - h_0(Z') = L_0(Z - Z'),$$

and (3) one has for any $\varepsilon > 0$ and $Z'' \in \mathcal{D}'$ where $0 \leq Z - Z'' \leq \varepsilon$ that

$$\frac{L_0}{1 + a\varepsilon}(Z - Z'') \leq h_0(Z) - h_0(Z''), \quad (17)$$

with $a \geq 0$.

Hence a denotes how 'smooth' the constant L_0 decays in a neighborhood of Z where the actual Lipschitz constraint is met (that is, a smaller a indicates higher smoothness) (see Figure 14). In particular, a value $a \rightarrow 0$ (arbitrarily small) implies that the function h_0 is linear with slope L_0 . Note that this definition does not require that the function $\frac{\partial h_0(Z)}{\partial Z}$ exists for any $Z \in \mathbb{R}$. This definition implies the inequality

$$\frac{1}{L_0} = \frac{Z - Z'}{h_0(Z) - h_0(Z')} \leq \frac{Z - Z''}{h_0(Z) - h_0(Z'')} \leq \frac{Z - Z'}{h_0(Z) - h_0(Z')} + \frac{a\varepsilon}{L_0}. \quad (18)$$

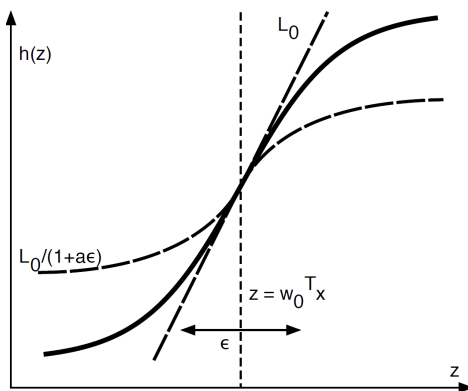


Figure 14: *Illustration of a function $h : \mathbb{R} \rightarrow \mathbb{R}$ (solid curved line) which is (L_0, a) -Lipschitz monotone according to Definition 6. The Lipschitz inequality is tight around Z indicated by the dotted line, while in the ϵ -neighborhood of Z the actual Lipschitz constant decays as slowly as $\frac{L_0}{1+a\epsilon}$.*

The first inequality holds due to the Lipschitz smoothness of h_0 . The second inequality follows from (17).

We now state that (w_0, h_0) can be recovered successfully ('identified') if \mathcal{D} and h_0 were such that Definition 5 and 6 hold. We consider the MINLIP estimator based on an ϵ -non-degenerate set \mathcal{D} which is defined as

$$w_\epsilon = \arg \max_{\|w\|_2=1} \inf_{(X,Y), (X',Y') \in \mathcal{D}: Y > Y'} \frac{w^T (X - X')}{Y - Y'},$$

or equivalently (up to a scaling)

$$\begin{aligned} w_\epsilon \propto \arg \min_w \frac{1}{2} w^T w \\ \text{s.t. } (Y - Y') \leq w^T (X - X') \quad \forall (X, Y), (X', Y') \in \mathcal{D} : Y > Y'. \end{aligned} \quad (19)$$

Specifically $w_\epsilon = \frac{\hat{w}}{\|\hat{w}\|_2}$ where \hat{w} is the optimizer of (19). If \mathcal{D} contains a finite number of elements this problem can be solved efficiently as a convex Quadratic Program (QP) using standard solvers. This estimator would return the desired result $w_\epsilon = w_0$ if enough observations were given. This is stated more formally as follows.

Lemma 2 (Identifiability) *Let $\epsilon > 0$ be any (arbitrarily small) constant. Given a model (h_0, w_0) governing the observations in \mathcal{D} . Assume that (i) the set \mathcal{D} is $(a\epsilon)$ -non-degenerate as in Definition 5; (ii) the function h_0 is (L_0, a) -Lipschitz monotone on the set $\mathcal{D}' = \{Z = w_0^T X \in \mathbb{R} : (X, Y) \in \mathcal{D}\}$, as in Definition 6. Then one has for all $w \in \mathbb{R}^d$ where $\|w\|_2 = 1$ that*

$$\frac{1}{L_0} \leq \inf_{(X,Y), (X',Y') \in \mathcal{D}: Y > Y'} \frac{w^T (X - X')}{Y - Y'},$$

with equality if $w = w_0$.

Proof Let $(X, Y), (X', Y') \in \mathcal{D}$ be such that $Y \neq Y'$ and the Lipschitz constant is achieved, or

$$\frac{1}{L_0} = \frac{w_0^T(X - X')}{h_0(w_0^T X) - h_0(w_0^T X')} = \frac{w_0^T(X - X')}{Y - Y'}. \quad (20)$$

Such an observation exists assuming that h_0 is (L_0, a) -Lipschitz on \mathcal{D}' . We prove that for a bad estimation w of w_0 ($w^T w_0 < 1 - \varepsilon$), one can always find an observation $(X^*, Y^*) \in \mathcal{D}$ such that $\frac{w^T(X - X^*)}{Y - Y^*}$ is strictly lower than $\frac{1}{L_0}$. This implies that when w deviates a fraction ε from w_0 , the objective in (19) can never achieve the maximum value as would be the case when $w = w_0$. This implies consistency of the MINLIP estimator.

At first, by the (L_0, a) -Lipschitz condition on h_0 , one has for all $(X'', Y'') \in \mathcal{D}$ where $|w_0^T(X - X'')| \leq \varepsilon$ and $Y \neq Y''$ that (as in inequality (18)),

$$\frac{w_0^T(X - X')}{Y - Y'} \geq \frac{w_0^T(X - X'')}{Y - Y''} - \frac{a\varepsilon}{L_0}. \quad (21)$$

According to the Cauchy-Schwarz' inequality, the condition $|w_0^T(X - X'')| \leq \varepsilon$ is fulfilled for $\|X - X''\|_2 \leq \varepsilon$. Secondly, for any $w \in \mathbb{R}^d$ with $\|w\|_2 = 1$ and $w_0^T w < 1 - a\varepsilon$, one has by the orthogonal decomposition of a vector that

$$\begin{aligned} w_0 - w &= \frac{w_0 w_0^T}{\|w_0\|_2^2} (w_0 - w) + v \\ &= w_0 (w_0^T w_0 - w_0^T w) + v \\ &= w_0 a\varepsilon^+ + v, \end{aligned}$$

with v the orthogonal complement of the projection of $w_0 - w$ on w_0 and $\varepsilon^+ > \varepsilon$. It follows for any $(X''', Y''') \in \mathcal{D}$ where $Y \neq Y'''$ that

$$\frac{(w_0 - w)^T(X - X''')}{Y - Y'''} = \frac{a\varepsilon^+ w_0^T(X - X''')}{Y - Y'''} + \frac{v^T(X - X''')}{Y - Y'''}.$$

Hence by assumption of the set \mathcal{D} being $(a\varepsilon)$ -non-degenerate, there exists for any $w \in \mathbb{R}^d$ (and thus for any $v \in \mathbb{R}^d$) an observation $(X^*, Y^*) \in \mathcal{D}$ with $\|X - X^*\|_2 \leq a\varepsilon$, $Y \neq Y^*$ such that

$$\frac{(w_0 - w)^T(X - X^*)}{Y - Y^*} = a\varepsilon^+ \frac{w_0^T(X - X^*)}{Y - Y^*} > a\varepsilon \frac{w_0^T(X - X^*)}{Y - Y^*} \geq \frac{a\varepsilon}{L_0}. \quad (22)$$

From (21) and (22) it then follows that

$$\frac{1}{L_0} = \frac{w_0^T(X - X')}{Y - Y'} \geq \frac{w_0^T(X - X^*)}{Y - Y^*} - \frac{a\varepsilon}{L_0} > \frac{w^T(X - X^*)}{Y - Y^*}.$$

Hence, for all $w \in \mathbb{R}^d$ for which $\|w\|_2 = 1$ and $w_0^T w < 1 - a\varepsilon$, there are two observations $(X, Y), (X^*, Y^*) \in \mathcal{D}$ such that

$$\frac{1}{L_0} > \frac{w^T(X - X^*)}{Y - Y^*},$$

proving the result. Equality as in (20) is reached for $w = w_0$. ■

Appendix B. MINLIP for Ranking Problems

The formal derivation of the MINLIP method is given in this Appendix. We start with the problem formulation as denoted in Equation (8):

$$\begin{aligned} \min_{w, \varepsilon} \quad & \frac{1}{2} w^T w + \gamma \|\varepsilon\|_1 \\ \text{s.t.} \quad & \mathbf{D}(\Phi w + \varepsilon) \geq \mathbf{D}\mathbf{Y}, \end{aligned}$$

with $\Phi = [\varphi(X_1), \dots, \varphi(X_n)]^T$. Take $\varepsilon = e^+ + e^-$ and suppose $e^+ \geq 0$ and $e^- \geq 0$. The problem can then be formulated as

$$\begin{aligned} \min_{w, e^+, e^-} \quad & \frac{1}{2} w^T w + \gamma \mathbf{1}_n^T (e^+ + e^-) \\ \text{s.t.} \quad & \begin{cases} \mathbf{D}(\Phi w + (e^+ - e^-)) \geq \mathbf{D}\mathbf{Y}, \\ e^+ \geq 0, \\ e^- \geq 0. \end{cases} \end{aligned}$$

The Lagrangian becomes

$$\mathcal{L}(w, e^+, e^-; \alpha, \beta_+, \beta_-) = \frac{1}{2} w^T w + \gamma \mathbf{1}_n^T (e^+ + e^-) - \beta_+^T e^+ - \beta_-^T e^- - \alpha^T \mathbf{D}(\Phi w + e^+ - e^- - \mathbf{Y}),$$

with Lagrange multipliers $\alpha, \beta^+, \beta^- \geq 0$. The conditions for optimality (Karush-Kuhn-Tucker (KKT) conditions) become

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial w} = 0 \rightarrow w = (\Phi)^T \mathbf{D}^T \alpha \\ \frac{\partial \mathcal{L}}{\partial e^+} = 0 \rightarrow \gamma = \mathbf{D}^T \alpha + \beta^+ \\ \frac{\partial \mathcal{L}}{\partial e^-} = 0 \rightarrow \gamma = -\mathbf{D}^T \alpha + \beta^- \\ \text{diag}(\alpha) \mathbf{D}(\Phi w + e^+ - e^- - \mathbf{Y}) = 0 \\ \text{diag}(\beta^+) e^+ = 0 \\ \text{diag}(\beta^-) e^- = 0 \\ \alpha \geq 0 \\ \beta^+ \geq 0 \\ \beta^- \geq 0, \end{cases} \quad (23)$$

where $\text{diag}(a)$ indicates a diagonal matrix with the elements of the vector a on the main diagonal. Now from Slater's condition one could exchange $\min_{w, e^+, e^-} \max_{\alpha}$ with $\max_{\alpha} \min_{w, e^+, e^-}$. Solving for w, e^+ and e^- gives the dual problem

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T \mathbf{D} \mathbf{K} \mathbf{D}^T \alpha - \alpha^T \mathbf{D}\mathbf{Y} \\ \text{s.t.} \quad & \begin{cases} -\gamma \mathbf{1}_n \leq \mathbf{D}^T \alpha \leq \gamma \mathbf{1}_n \\ \alpha \geq 0_{n-1}, \end{cases} \end{aligned}$$

and from the first condition of (23) and the model specification $u(X) = w^T \varphi(X)$ one could write the solution for a new point X^* as

$$\hat{u}(X^*) = \mathbf{K}_n^* \mathbf{D}^T \hat{\alpha},$$

with $\mathbf{K}_n^* \in \mathbb{R}^n$ and $\mathbf{K}_n^* = [K(X^*, X_1) \dots K(X^*, X_n)]^T$.

Appendix C. MINLIP for Ordinal Regression

In this appendix the derivation of the MINLIP method for ordinal regression is exposed. In the ordinal regression case unknown thresholds v are introduced corresponding to an outcome intermediate between two successive outcome levels. The model is built by indicating that the difference between the utility of a certain observation X_i and the largest threshold lower than the outcome of that observation Y_i should be larger than the difference between Y_i and the outcome corresponding to the before mentioned threshold. Analogously, the difference between the smallest threshold higher than Y_i should be larger than the difference between the outcome corresponding to that threshold and Y_i . As an extra constraint we impose that successive threshold are increasing values of the utility function. More formally the problem is formulated as in Equation (11), now using the kernel based version:

$$\begin{aligned} \min_{w, e, e^*, v} \quad & \|w\|_2 + \gamma \mathbf{1}_n^T (e + e^*) \\ \text{s.t.} \quad & \begin{cases} \Phi w - \mathbf{Q}v + e \geq \mathbf{Y} - \mathbf{Q}\mathbf{B} \\ -\Phi w + \mathbf{Q}^*v + e^* \geq -\mathbf{Y} + \mathbf{Q}^*\mathbf{B} \\ e \geq 0 \\ e^* \geq 0 \\ \mathbf{M}v \leq 0. \end{cases} \end{aligned}$$

As in Appendix B we build the Lagrangian

$$\begin{aligned} \mathcal{L}(w, e^+, e^-; \alpha, \beta, \eta, \eta^*, v) = & \frac{1}{2} w^T w + \gamma \mathbf{1}_n^T (e + e^*) - \alpha^T (\Phi w - \mathbf{Q}v + e - \mathbf{Y} + \mathbf{Q}\mathbf{B}) \\ & - \beta^T (-\Phi w + \mathbf{Q}^*v + e^* + \mathbf{Y} - \mathbf{Q}^*\mathbf{B}) - \eta^T e - \eta^{*T} e^* \\ & + v^T \mathbf{M}v, \end{aligned}$$

and derive the set of optimality conditions

$$\left\{ \begin{array}{l} \frac{\partial \mathcal{L}}{\partial w} = 0 \rightarrow w = \Phi^T (\alpha - \beta) \\ \frac{\partial \mathcal{L}}{\partial e} = 0 \rightarrow \gamma \mathbf{1}_n = \alpha + \eta \\ \frac{\partial \mathcal{L}}{\partial e^*} = 0 \rightarrow \gamma = \beta + \eta^* \\ \frac{\partial \mathcal{L}}{\partial v} = 0 \rightarrow \alpha^T \mathbf{Q} - \beta^T \mathbf{Q}^* + v^T \mathbf{M} = 0 \\ \text{diag}(\alpha) D(\Phi w - \mathbf{Q}v + e - \mathbf{Y} + \mathbf{Q}\mathbf{B}) = 0 \\ \text{diag}(\beta) (-\Phi w + \mathbf{Q}^*v + e^* + \mathbf{Y} - \mathbf{Q}^*\mathbf{B}) = 0 \\ \text{diag}(\eta) e = 0 \\ \text{diag}(\eta^*) e^* = 0 \\ \text{diag}(v) \mathbf{M}v = 0 \\ \alpha \geq 0 \\ \beta \geq 0 \\ \eta \geq 0 \\ \eta^* \geq 0 \\ v \geq 0. \end{array} \right.$$

The dual problem formulation is than found as

$$\begin{aligned} \min_{\alpha, \beta} \quad & \frac{1}{2} \alpha^T \mathbf{K} \alpha + \frac{1}{2} \beta^T \mathbf{K} \beta - \alpha^T \mathbf{K} \beta - \alpha^T (\mathbf{Y} - \mathbf{B}^T \mathbf{Q}) + \beta^T (\mathbf{Y} - \mathbf{B}^T \mathbf{Q}^*) \\ \text{s.t.} \quad & \begin{cases} 0_n \leq \alpha \leq \gamma \mathbf{1}_n \\ 0_n \leq \beta \leq \gamma \mathbf{1}_n \\ 0_{k-2} \leq \mathbf{v} \\ \mathbf{Q}^T \alpha - \mathbf{Q}^{*T} \beta + \mathbf{M}^T \mathbf{v} = 0_{k-1}. \end{cases} \end{aligned}$$

References

- N. Ailon and M. Mohri. An efficient reduction of ranking to classification. In *COLT*, pages 87–98. Omnipress, 2008.
- P.K. Andersen, O. Borgan, R.D. Gill, and N. Leiding. *Statistical Models based on Counting Processes*. Springer-Verlag, New York, 1993.
- E. Bair and R. Tibshirani. Semi-supervised methods to predict patient survival from gene expression data. *PLoS Biology*, 2(4):511–522, April 2004.
- R. Bair, T. Hastie, P. Debnath, and R. Tibshirani. Prediction by supervised principal components. *Journal of the American Statistical Association*, 101:119–137, 2006.
- E. Biganzoli, P. Boracchi, L. Mariani, and E. Marubini. Feedforward neural networks for the analysis of censored survival data: a partial logistic regression approach. *Statistics in Medicine*, 17(10):1169–1186, 1998.
- H.M.M. Bøvelstad, S. Nygård, H.L.L. Størvold, M. Aldrin, O. Borgan, A. Frigessi, and O.C.C. Lingjærde. Predicting survival from microarray data - a comparative study. *Bioinformatics*, 23(16):2080–2087, 2007.
- N. Breslow. Covariance analysis of censored survival data. *Biometrics*, 30(1):89–99, 1974.
- S.C. Cheng, L.J. Wei, and Z. Ying. Predicting survival probabilities with semiparametric transformation models. *Journal of the American Statistical Association*, 92(437):227–235, 1997.
- W. Chu and Z. Ghahramani. Gaussian processes for ordinal regression. *Journal of Machine Learning Research*, 6:1019–1041, 2005.
- W. Chu and S.S. Keerthi. New approaches to support vector ordinal regression. In *In ICML 2005: Proceedings of the 22nd international conference on Machine Learning*, pages 145–152, 2005.
- S. Cléménçon and N. Vayatis. Ranking the best instances. *Journal of Machine Learning Research*, 8:2671–2699, 2007.
- S. Cléménçon, G. Lugosi, and N. Vayatis. Ranking and scoring using empirical risk minimization. In *Proceedings of the 18th Annual Conference in Learning Theory (COLT)*, Bertinoro (Italy), June 2005.
- D.R. Cox. Regression models and life-tables (with discussion). *Journal of the Royal Statistical Society, Series B*, 34(2):187–220, 1972.

- D.M. Dabrowska and K.A. Doksum. Partial likelihood in transformation models with censored data. *Scandinavian Journal of Statistics*, 15(1):1–23, 1988.
- P.H. Eilers and B.D. Marx. Flexible smoothing with B-splines and penalties. *Statistical Science*, 11:89–121, 1996.
- R.C. Elandt-Johnson and N.L. Johnson. *Survival Models and Data Analysis*. John Wiley & Sons, Inc., 1980.
- D.M. Finkelstein. A proportional hazards model for interval-censored failure time data. *Biometrics*, 42:845–854, 1986.
- B. Fisher, M. Bauer, L. Wickerham, C.K. Redmong, and E.R. Fisher. Relation of number of positive axillary nodes to the prognosis of patients with primary breast cancer. An NSABP update. *Cancer*, 52(9):1551–1557, 1983.
- Y. Freund, R. Iyer, R.E. Schapire, and Y. Singer. An efficient boosting algorithm for combining preferences. *Journal of Machine Learning Research*, 4(6):933–969, 2004.
- J. Fürnkranz and E. Hüllermeier. Pairwise preference learning and ranking. In *Proceedings of the European Conference on Machine Learning 2003, Cavtat-Dubrovnik*, pages 145–156. Springer-Verlag, 2003.
- F. Harrell. *Regression Modeling Strategies. With Applications to Linear Models, Logistic Regression, and Survival Analysis*. Springer, 2001.
- F. Harrell, K. Klee, R. Califf, D. Pryor, and R. Rosati. Regression modeling strategies for improved prognostic prediction. *Statistics in Medicine*, 3(2):143–152, 1984.
- T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. Springer-Verlag, 2001.
- P. Heagerty, T. Lumley, and M. Pepe. Time-dependent ROC curves for censored survival data and a diagnostic marker. *Biometrics*, 56(2):337–344, 2000.
- R. Herbrich, T. Graepel, P. Bollmann-Sdorra, and K. Obermayer. Learning preference relations for information retrieval. In *Proceedings of the Fifteenth Conference of the American Association for Artificial Intelligence*, pages 1–4, 1998.
- R. Herbrich, T. Graepel, and K. Obermayer. Large margin rank boundaries for ordinal regression. In Smola, Bartlett, Schoelkopf, and Schuurmans, editors, *Advances in Large Margin Classifiers*, 2000.
- C.M. Hurvich, J.S. Simonoff, and C.L. Tsai. Smoothing parameter selection in nonparametric regression using an improved akaike information criterion. *Journal of the Royal Statistical Society: Series B*, 60:371–293, 1998.
- J.D. Kalbfleisch and R.L. Prentice. *The Statistical Analysis of Failure Time Data*. Wiley series in probability and statistics. 2002.

- M.W. Kattan, H. Ishida, P.T. Scardino, and J.R. Beck. Applying a neural network to prostate cancer survival data. In N. Lavrac, R. Keravnou, and B. Zupan, editors, *Intelligent Data Analysis in Medicine and Pharmacology*, pages 295–306. Kluwer, Boston, 1997.
- R. Koenker and O. Geling. Reappraising medfly longevity: A quantile regression survival analysis. *Journal of the American Statistical Association*, 96:458–468, 2001.
- P.J. Lamy, P. Pujol, S. Thezenas, A. Kramar, P. Rouanet, F. Guilleux, and J. Grenier. Progesterone receptor quantification as a strong prognostic determinant in postmenopausal breast cancer women under tamoxifen therapy. *Breast cancer reasearch and treatment*, 76(1):65–71, 2002.
- D.J.C. MacKay. A practical Bayesian framework for backpropagation networks. *Neural Computing*, 4(3):448–472, 1992.
- H. Martens and T. Næs. *Multivariate Calibration*. New York: John Wiley & Sons Inc., 1989.
- R.G. Miller. *Survival Anlysis*. John Wiley & Sons, 1981.
- S. Nygård, O. Borgan, O. Lingjærde, and H. Størvold. Partial least squares Cox regression for genome-wide data. *Lifetime Data Analysis*, 14(2):179–195, 2008.
- K. Pelckmans, M. Espinoza, J. De Brabanter, J.A.K. Suykens, and B. De Moor. Primal-dual monotone kernel regression. *Neural Processing Letters*, 22(2):171–182, 2005a.
- K. Pelckmans, I. Goethals, J. De Brabanter, J.A.K. Suykens, and B. De Moor. *Componentwise Least Squares Support Vector Machines*, chapter in Support Vector Machines: Theory and Applications, pages 77–98. (L. Wang, ed.), Springer, 2005b.
- R. Peto. Discussion of paper by D.R. Cox. *Journal of the Royal Statistical Society, Series B*, 34: 205–207, 1972.
- M.F. Pichon, C. Pallud, M. Brunet, and E. Milgrom. Relationship of presence of progesterone receptors to prognosis in early breast cancer. *Cancer Research*, 40:3357–3360, 1980.
- A. Rosenwald, G. Wright, W.C. Chan, J.M. Connors, E. Campo, R.I. Fisher, R.D. Gascoyne, H.K. Muller-Hermelink, R.B. Smeland, J.M. Giltneane, E.M. Hurt, H. Zhao, L. Averett, and L. Yang. The use of molecular profiling to predict survival after chemotherapy for diffuse large-B-cell lymphoma. *The New England Journal of Medicine*, 346(25):1937–1947, 2002.
- G.A. Satten. Rank-based inference in the proportional hazards model for interval censored data. *Biometrika*, 83:355–370, 1996.
- M. Schumacher, G. Basert, H. Bojar, K. Huebner, M. Olschewski, W. Sauerbrei, C. Schmoor, C. Beyerle, R.L.A. Neumann, and H.F. Rauschecker. Randomized 2 x 2 trial evaluating hormonal treatment and the duration of chemotherapy in node-positive breast cancer patients. *Journal of Clinical Oncology*, 12, 1994.
- J. Shawe-Taylor and N. Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.

- T. Sørli, R. Tibshirani, J. Parker, T. Hastie, J.S. Marron, A. Nobel, S. Deng, H. Johnsen, R. Pesich, S. Geisler, J. Demeter, C.M. Perou, P.E. Lønning, P.O. Brown, A. Børresen-Dale, and D. Botstein. Repeated observation of breast tumor subtypes in independent gene expression data sets. *Proceedings of the National Academy of Sciences of the United States of America*, 100(14):8418–8423, 2003.
- J.A.K. Suykens, T. Van Gestel, J. De Brabanter, B. De Moor, and J. Vandewalle. *Least Squares Support Vector Machines*. World Scientific, Singapore, 2002.
- T.M. Therneau and P.M. Grambsch. *Modeling Survival Data: Extending the Cox Model*. Springer, 2 edition, 2000.
- R. Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society, Series B*, 58(1):267–288, 1996.
- V. Van Belle, K. Pelckmans, J.A.K. Suykens, and S. Van Huffel. Support vector machines for survival analysis. In *Proceedings of the Third International Conference on Computational Intelligence in Medicine and Healthcare (CIMED2007)*, pages 1–8, Plymouth (UK), July 2007.
- V. Van Belle, K. Pelckmans, J.A.K. Suykens, and S. Van Huffel. Survival SVM: a practical scalable algorithm. In *Proceedings of the 16th European Symposium on Artificial Neural Networks (ESANN2008)*, pages 89–94, Bruges (Belgium), April 2008.
- V. Van Belle, K. Pelckmans, J.A.K. Suykens, and S. Van Huffel. MINLIP: Efficient learning of transformation models. In *Proceedings of the International Conference on Artificial Neural Networks (ICANN2009)*, pages 60–69, Limassol (Cyprus), September 2009.
- M.J. van de Vijver, L.J. van’t Veer, and H. Dai. A gene-expression signature as a predictor of survival in breast cancer. *The New England Journal of Medicine*, 347(25):1999–2009, 2002.
- H.C. van Houwelingen, T. Bruinsma, A.A.M. Hart, L.J. van’t Veer, and L.F.A. Wessels. Cross-validated cox regression on microarray gene expression data. *Statistics in Medicine*, 25(18):3201–3216, 2006.
- V. Vapnik. *Statistical Learning Theory*. Wiley and Sons, 1998.
- C. Verschraegen, C. Vihn-Hung, G. Cserni, R. Gordon, M.E. Royce, G. Vlastos, P. Tai, and G. Storme. Modeling the effect of tumor size in early breast cancer. *Annals of Surgery*, 241(2):309–318, 2005.