# Ordinal Risk-Group Classification

Yizhar Toren

Tel Aviv University, Tel Aviv, Israel yizhar.toren@math.tau.ac.il

June 13, 2018

### Abstract

Most classification methods provide either a prediction of class membership or an assessment of class membership probability. In the case of two-group classification the predicted probability can be described as "risk" of belonging to a "special" class . When the required output is a set of ordinal-risk groups, a discretization of the continuous risk prediction is achieved by two common methods: by constructing a set of models that describe the conditional risk function at specific points (quantile regression) or by dividing the output of an "optimal" classification model into adjacent intervals that correspond to the desired risk groups. By defining a new error measure for the distribution of risk onto intervals we are able to identify lower bounds on the accuracy of these methods, showing sub-optimality both in their distribution of risk and in the efficiency of their resulting partition into intervals. By adding a new form of constraint to the existing maximum likelihood optimization framework and by introducing a penalty function to avoid degenerate solutions, we show how existing methods can be augmented to solve the ordinal risk-group classification problem. We implement our method for generalized linear models (GLM) and show a numeric example using Gaussian logistic regression as a reference.

## 1 Introduction

The classical problem of discriminating between two classes of observations based on a given dataset has been widely discussed in the statistical literature. When only two classes are involved, the question of discrimination is reduced to whether or not a given observation is a member of a "special" class (where the other class is the default state, for example sick vs. healthy). Some classification methods, such as Fisher's linear discriminant analysis (LDA), make a decisive prediction of class membership while minimizing error in some sense, typically the misclassification rate. Other methods, such as logistic regression, provide an estimate of the exact conditional probability of belonging to the "special" class given a set of predictor variables. Throughout this paper we shall refer to this conditional probability as "*conditional risk*" or simply "*risk*", although sometimes belonging to the special class might actually have a very positive context (e.g. success).

There are two ways to estimate the conditional risk function: parametric and non-parametric. Parametric methods primarily include logit/probit models (Martin 1977 [18], Ohlsen 1980 [21]) and linear models (Amemiya 1981 [1], Maddala 1986 [14] and Amemiya 1985 [2]). Powell (1994 [22]) has a review of non-parametric estimators. For a comparison of these approaches and complete review see Elliott and Lieli (2006) [7] and more recently Green and Hensher (2010) [10].

The estimation of the exact structure of the conditional risk function comes in handy when we wish to make distinctions between observations that are finer than simply class membership. However, in realistic scenarios acting upon such estimations alone may prove to be difficult. Assessments on a scale of 1:100 (as percentages) or finer assessments have little practical use, primarily since the possible actions resulting from such information are usually few. For such cases an ordinal output is required. It is important to note that this problem is not equivalent to multi-group classification in two ways: first, our groups are ordinal by nature and relate to the underlying risk; second, the assignment into groups is not given a-priori and greatly depends on the selection of model, model parameters and the borders of the intervals assigned to each risk group.

There are two common approaches to creating an ordered set of risk groups to match a finite set of escalating actions. The first approach is to create multiple models describing the behaviour of the conditional risk function at specific points (also known as "quantile regression"); the second approach is to divide post-hoc the continuous risk estimation of a known model into intervals.

The first approach attempts to construct separate models that describe the behaviour of the conditional risk function at specific levels of risk. In linear models this approach is known as *quantile regression* (Koenker & Bassett 1978 [13]). Manski ([15], [16], [17]) implemented this notion to binary response models (the equivalent of two-group classification) naming it "Maximum Score Estimation". In a series of papers he shows the existence, uniqueness and optimal properties of the estimators and follows by showing their stable asymptotic properties. The primary justification for using this approach is methodological: it demands that we specify in advance the levels of risk that are of interest to us (a vector $q$ of quantiles), and then constructs a series of models that describe conditional risk at these quantiles. However, as we shall demonstrate in section 3.1, using risk-quantiles (or conditional probability over left-unbounded and overlapping intervals) is not relevant to our definition of the problem and even the term "conditional quantiles" is in itself misleading.

In the second, more "practical" approach, the continuous output of an existing optimal risk model (logit, linear or non-parametric) is divided into intervals, thus translating the prediction of risk (usually continuous in nature) into large "bins of risk" - i.e "low"/"medium"/"high" or "mild"/"moderate"/"severe" (depending on context). The final result of this discretization process is a set of ordinal risk groups based on the continuous prediction of conditional risk. The primary drawback of this approach is that the selection of the classification model and its parameters is not performed in light of the final set of desired risk groups. Instead, an "optimal model" (in some sense) is constructed first, and the partition into discrete groups is performed post-hoc.

The primary objective of this paper is to combine the idea of pre-set levels of risk over adjacent intervals (rather than risk quantiles) into a standard classification framework. Instead of constructing multiple models, we offer a process that optimizes a single risk estimation model (or "score") paired with a matching set of breakpoints that partition the model's output into ordinal risk groups. To that end we define a new measure of accuracy - *Interval Risk Deviation* (IRD) - which describes a model's ability to distribute risk correctly into intervals given a pre-set vector $r$ of risk levels. We show how this new measure of error can be integrated into existing classification frameworks (specifically the maximum likelihood framework) by adding a constraint to the existing optimization problem. In addition, we address the more practical problem of effectively selecting breakpoints by introducing a penalty function to the modified optimization scheme.

The remainder of this paper is organized as follows. Section 2 defines risk groups and a measure of error (IRD) that will be necessary for optimality. Section 3 demonstrates the problems of using existing approaches. Section 4 formulates a new optimization problem that will provide accurate, optimal and non-degenerate solutions, and section 5 provides a case study where the new framework is applied to logistic regression and presents an example.

# 2 Definitions

Let $r \in [0,1]^T$ be an ordered vector of *risk levels* ($0 \le r_1 < r_2 < \ldots < r_T \le 1$), let $X = (X_1, \ldots, X_P)$ be a continuous $P$-dimensional random vector and let $Y \in \{0,1\}$ be a Bernoulli random variable representing class membership. An *Ordinal Risk-Group Score* (ORGS) for a pre-set risk vector $r$ is a couplet $(\Psi, \tau)$ where $\Psi : \mathbb{R}^P \to \mathbb{R}$ is a continuous (possibly not normalized) risk predictor, which summarizes the attributes of $X$ into a single number (a score), and $\tau \in \mathbb{R}^{T-1}$ is a complete partition of $\mathbb{R}$ into $T$ distinct and adjacent intervals ($-\infty = \tau_0 < \tau_1 < \tau_2 < \ldots < \tau_{T-1} < \tau_T = \infty$). The couplet $(\Psi, \tau)$ classifies observations into risk groups by the following equivalence: An observed vector $X$ belongs to the $i$'th risk group if and only if $\Psi(X) \in (\tau_{i-1}, \tau_i]$ (the intervals are right-side open to avoid ambiguities). The actual conditional risk level of the $i$'th risk group defined by a couplet $(\Psi, \tau)$ is:

$$R_i(\Psi, \tau) = P(Y = 1 \mid \Psi(X) \in (\tau_{i-1}, \tau_i]) \tag{1}$$

It is worth noting that score-based classification methods for two classes can be described as a special of $T = 2$ (two risk groups). Such methods look for a single breakpoint $\tau \in \mathbb{R}$, and the two resulting intervals $(-\infty, \tau], (\tau, \infty)$ become an absolute prediction of class membership: $\Psi(X) > \tau \Rightarrow X$ belongs to class 1. Other methods, designed to deal with more than one risk group, typically assign a single breakpoint to each risk group (see section 3.1), reflecting the idea that the assignment to risk group is based on *thresholds*: an observed $X$ is assigned to the $i$'th group if and only if $\Psi(X)$ crosses the $(i-1)$'th threshold ($\Psi(X) > \tau_{i-1}$) but does not cross the $i$'th threshold ($\Psi(X) \le \tau_i$).

Even from the latter definition, it becomes evident that the assignment to groups is in fact based on *adjacent intervals* $\{(\tau_{i-1}, \tau_i]\}_i^{T-1}$ (rather than on right-side open ended

intervals defined by thresholds) ans that any breakpoint we set affects the definition of two intervals (and hence two risk groups). Although further on in this paper we shall discuss separate breakpoints in relation to risk groups in order to demonstrate the key problem that arises from the use of adjacent intervals (section 3.2), the notion of assigning intervals *simultaneously* rather than separate breakpoints should remain clear throughout this paper.

We can now describe the accuracy of an ordinal risk score $(\Psi, \tau)$ in relation to a pre-set vector $r$ as the overall difference between the pre-defined risk levels of $r$ and the actual conditional risk levels $R(\Psi, \tau)$. We define an error measure for risk-group classification models which is a parallel of *misclassification rate* in standard classification methods. We name this measure *Interval Risk Deviation* (IRD):

$$\text{IRD}_r(\Psi, \tau) = \|R(\Psi, \tau) - r\| \tag{2}$$

On it's own, the very definition of IRD marks a new approach to the evaluation of ordinal risk scores. Having a predefined set of risk levels means that any risk score $(\Psi, \tau)$ we consider as a candidate must uphold $\text{IRD}_r(\Psi, \tau) = 0$ (or at the very least $\text{IRD}_r(\Psi, \tau) < \varepsilon$ for a predefined small $\varepsilon > 0$). This makes IRD = 0 a *necessary condition* for optimality. In the next two sections we demonstrate how the two existing approaches for creating ordinal risk scores do not necessarily fulfil this condition, either because of unsuitable definitions of optimality, as is the case with risk-quantile based methods, or by ignoring it altogether, as is the case with the 2-step approach.

# 3 Problems with Existing Scoring Methods

## 3.1 Risk-Quantiles (and why we can't use them)

When first presented with the problem of selecting an optimal model paired with a set of optimal breakpoints, our initial idea was to use quantile-oriented models. Such models have been extensively studied in econometrics, where they are commonly referred to as "ordered choice models" ([27], [10]). The most relevant model in that group is Manski's *maximum score estimation* which defines the optimization problem using a set of probabilities over *left-unbounded overlapping* intervals (or *rays*) in contrast to the definition of the problem over *adjacent, non-overlapping* intervals.

In order to better illustrate the differences between our definitions and Manski's quantile-oriented approach we must first describe quantile oriented models in our terms. First we replace the vector $r$ with a vector $q$ of "conditional quantiles", which are in fact the desired conditional probabilities over left-unbounded and overlapping intervals. Using Manski's adaptation of quantile regression [15] we can build a different set of model parameters for each quantile $q_i$ optimizing:

$$|P(Y = 1 \mid \Psi_i(X) \leq 0) - q_i| \longrightarrow \min_{\Psi_i} \tag{3}$$

4

It is easy to see how this approach can be slightly modified to match the original objective of finding a single model: by coercing the models $\Psi_i$ to be parallel we can create a "master model" $\Psi(X)$ and derive appropriate thresholds $\{\tau_i\}_{i=1}^{T-1}$ such that:

$$\Psi_i(X) \leq 0 \quad \Leftrightarrow \quad \Psi(X) \leq \tau_i$$

$$|P(Y = 1 \mid \Psi(X) \leq \tau_i) - q_i| \longrightarrow \min_{\Psi} \quad i \in \{1, \ldots T\} \tag{4}$$

Using (4) we can easily define $Q_i(\Psi, \tau) = P(Y = 1 \mid \Psi(X) \leq \tau_i)$ and the equivalent *Quantile Risk Deviation* $QRD_q(\Psi, \tau) = \|Q(\Psi, \tau) - q\|$, and look for a model with $QRD = 0$. However, while it is tempting to describe the vector $q$ as a vector of "*conditional quantiles*", the term is in itself misleading and should be avoided. Figure 1 demonstrates how even under relatively simple assumptions (a one dimensional Gaussian distribution with unequal conditional variances) the function $Q_i(\Psi, \tau) = P(Y = 1 \mid X \leq \tau_i)$ is not even monotone in $\tau_i$.
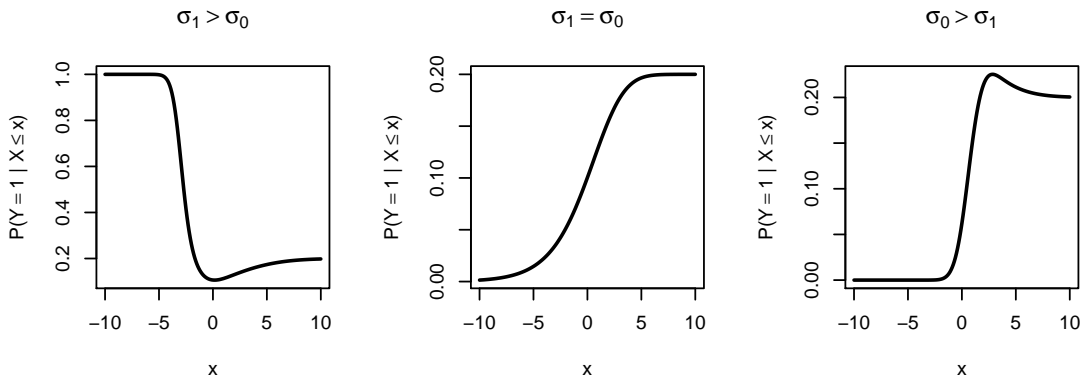


Figure 1: Different behaviour of conditional probability over left-unbounded intervals as a function of the threshold $x$ in the case of one-dimensional Gaussian distribution with $\mu_0 = -1$, $\mu_1 = 1$ and $P(Y = 1) = 0.2$. In the left panel $\sigma_1 = 4, \sigma_0 = 1$, in the middle panel $\sigma_1 = \sigma_0 = 2$ (homoscedastic case) and in the right panel $\sigma_1 = 1, \sigma_0 = 4$.

Even if we assume strict monotonicity of $P(Y = 1 \mid \Psi(X) \leq x)$, for example by assuming the strict monotone likelihood ratio property (SMLRP, for details see Appendix A) and thus giving the term "conditional quantiles" a meaningful sense, it would still be impossible to apply this approach to optimizing the distribution of risk over adjacent intervals. In order to use "risk-quantiles" to solve our problem we must first find an a-priori mechanism that will translate any given vector of desired conditional probabilities over adjacent intervals $r$ to the equivalent vector of desired conditional probabilities over left unbounded and overlapping intervals $q$.

However it is easy to show that such an a-priori translation is impossible. Using the *law of total probability* in its conditional form we can calculate for any given $R$ the

equivalent $Q^{(R)}$ (actual probabilities over left unbounded intervals):

$$
\begin{aligned}
Q_i^{(R)}(\Psi, \tau) &= P(Y = 1 \mid \Psi(X) \leq \tau_i) \\
&= \sum_{j \leq i} P(Y = 1 \mid \Psi(x) \in (\tau_{j-1}, \tau_j], \Psi(X) \leq \tau_i) P(\Psi(X) \in (\tau_{j-1}, \tau_j] \mid \Psi(X) \leq \tau_i) \\
&= \sum_{j \leq i} P(Y = 1 \mid \Psi(x) \in (\tau_{j-1}, \tau_j]) \frac{P(\Psi(X) \in (\tau_{j-1}, \tau_j], \Psi(X) \leq \tau_i)}{P(\Psi(X) \leq \tau_i)} \\
&= \frac{1}{P(\Psi(X) \leq \tau_i)} \sum_{j \leq i} R_j(\Psi, \tau) P(\Psi(X) \in (\tau_{j-1}, \tau_j])
\end{aligned}
$$

$$(5)$$

Or equivalently:

$$
R_i(\Psi, \tau) = \frac{P(\Psi(X) \leq \tau_i)}{P(\Psi(X) \in (\tau_{i-1}, \tau_i])} Q_i^{(R)}(\Psi, \tau) - \frac{P(\Psi(X) \leq \tau_{i-1})}{P(\Psi(X) \in (\tau_{i-1}, \tau_i])} Q_{i-1}^{(R)}(\Psi, \tau) \quad (6)
$$

The same process can be applied to the corresponding vector of risk quantiles $q^{(r)}$:

$$
\begin{aligned}
q_i^{(r)} &= \frac{\sum_{j<i} r_j \, P(\Psi(X) \in (\tau_{i-1}, \tau_i])}{P(\Psi(X) \leq \tau_i)} \\
r_i &= \frac{P(\Psi(X) \leq \tau_i)}{P(\Psi(X) \in (\tau_{i-1}, \tau_i])} q_i^{(r)} - \frac{P(\Psi(X) \leq \tau_{i-1})}{P(\Psi(X) \in (\tau_{i-1}, \tau_i])} q_{i-1}^{(r)}
\end{aligned}
$$

$$(7)$$

As a result for a fixed $(\Psi, \tau)$ we have:

$$
R_i(\Psi, \tau) = r_i \Leftrightarrow Q^{(R)}(\Psi, \tau) = q_i^{(r)}
$$

$$
\mathrm{IRD}_r(\Psi, \tau) = 0 \Leftrightarrow QRD_{q^{(r)}}(\Psi, \tau) = 0 \quad (8)
$$

The primary problem of using quantiles to define this problem stems from the relation between $r$ and the resulting $q_r$. By our own definitions the central aspect of the problem is the probability over adjacent intervals and not overlapping left-unbounded intervals. Therefore the optimization must be performed against a fixed, pre-defined vector $r$. If we wish to construct an analogous quantile-based optimization problem, we must first find the equivalent vector $q_r$ which defines quantile-based problem. However equation (7) shows that since the relation between $r$ and $q_r$ depends on the specific form of the optimal model $\Psi$, in order to construct $q_r$ we must first find the optimal model $\Psi$ for this problem (which is what we are looking for in the first place), or in other words the translation $r \leftrightarrow q$ is possible only once we have the optimal solution to the problem. Therefore building an analogous optimization problem over left-unbounded overlapping intervals can only be done *after* we have the optimal solution. Consequently we cannot use quantile-based models to construct an optimal model for the adjacent interval-based ordinal risk-group problem.

## 3.2   Lower bounds on Interval Risk Deviation

Another common practice when building scores for risk groups is to build a model $\Psi$ that is optimal in some sense (e.g. maximizing likelihood or minimizing overall miss-classification rate) and then partition the range of $\Psi(X)$ into adjacent intervals the define risk groups. In this section we demonstrate how, under relatively simple assumptions, using this approach with existing classification models is not optimal for more than two risk groups.

Using Bayes theorem we can represent $R$ as:

$$R_i(\Psi, \tau) = P(Y = 1 \mid \Psi(X) \in (\tau_{i-1}, \tau_i]) = P(Y = 1)\frac{P(\Psi(X) \in (\tau_{i-1}, \tau_i] \mid Y = 1)}{P(\Psi(X) \in (\tau_{i-1}, \tau_i])}$$

We assume that $(X, Y, \Psi)$ satisfies the Strict Monotone Likelihood Ratio Property (SMLRP, see appendix A for exact definition and details) and that the marginal densities $f_{X|Y=k}$ ($k = 0, 1$) are continuous, strictly positive and finite. By continuity and finiteness we can describe the behaviour of $R_i(\Psi, \tau)$ for infinitely short intervals $(\tau_i \to \tau_{i-1})$:

$$\lim_{\tau_i \to \tau_{i-1}} R_i(\Psi, \tau) = \lim_{\tau_i \to \tau_{i-1}} \frac{P(Y = 1)\, P(\Psi(X) \in (\tau_{i-1}, \tau_i] \mid Y = 1)}{P(\Psi(X) \in (\tau_{i-1}, \tau_i])} =$$

$$= P(Y = 1)\frac{\lim_{\tau_i \to \tau_{i-1}} \frac{P(\Psi(X) \in (\tau_i, \tau_{i-1}]|Y=1)}{\tau_i - \tau_{i-1}}}{\lim_{\tau_i \to \tau_{i-1}} \frac{P(\Psi(X) \in (\tau_i, \tau_{i-1}])}{\tau_i - \tau_{i-1}}} = P(Y = 1)\frac{f_{\Psi(X)|Y=1}(\tau_{i-1})}{f_{\Psi(X)}(\tau_{i-1})} \tag{9}$$

where $f$ is the appropriate density function and the limit is from the right-hand side. Similarly for any $z \in (\tau_{i-1}, \tau_i]$,

$$\lim_{\tau_i \to z} \lim_{\tau_{i-1} \to z} R_i(\Psi, \tau) = \lim_{\tau_{i-1} \to z} \lim_{\tau_i \to z} R_i(\Psi, \tau) = P(Y = 1)\frac{f_{\Psi(X)|Y=1}(z)}{f_{\Psi(X)}(z)} \tag{10}$$

Although we have stressed the importance of simultaneity when assigning intervals to risk groups, in order to understand the implications of (10) on optimal model selection we must look at the problem from a different perspective. First we fix $\Psi$ and assume that a given partition $\tau$ supports a perfect distribution of conditional risk up to the $(i-1)$'th group, meaning that $R_j(\Psi, \tau) = r_j$ for all $j < i$. Under these conditions, combined with our previous assumptions of continuous, strictly positive conditional densities and SMLRP, we can explicitly show that not all values of $r_i$ are exactly achievable without introducing some IRD: by theorem A.1 $R_i(\Psi, \tau)$ is strictly increasing in $\tau_i$ and therefore we can explicitly define a feasibility criterion:

$$P(Y = 1)\frac{f_{\Psi(X)|Y=1}(\tau_{i-1})}{f_{\Psi(X)}(\tau_{i-1})} < r_i \tag{11}$$

Using continuity (which enables us to divide by $P(Y = 1)f_{\Psi(X)|Y=1}(\tau_{i-1})$) we can transform (11) into a condition on the likelihood ratio $\Lambda$:

$$\Lambda_\Psi(\tau_{i-1}) = \frac{f_{\Psi(X)|Y=1}(\tau_{i-1})}{f_{\Psi(X)|Y=0}(\tau_{i-1})} < \frac{1 - P(Y = 1)}{P(Y = 1)}\frac{r_i}{1 - r_i} \tag{12}$$

If $\tau$ does not meet the feasibility criterion (11), then by (9) and strict monotonicity of $R$ any selection of $\tau_i > \tau_{i-1}$ will have $R_i(\Psi, \tau) > r_i$ even if we set the interval $(\tau_{i-1}, \tau_i]$ to be arbitrarily small. The inevitable result that, for the our fixed model $\Psi$, *any* choice of $\tau$ will have $\mathrm{IRD}_r(\Psi, \tau) > 0$.

It is important to note that the set of $T - 1$ inequalities defined by (11), (12) are necessary yet not sufficient conditions for IRD=0. Assume that we have a solution $(\Psi, \tau)$ which satisfies $\mathrm{IRD}_r(\Psi, \tau) = 0$. Under SMLRP we have $x_2 > x_1 \Rightarrow \Lambda_\Psi(x_2) > \Lambda_\Psi(x_1)$. Our counter example $(\Psi, \tilde{\tau})$ satisfies $\tilde{\tau}_1 < \tau_1$ and $\forall i > 1 : \tilde{\tau}_i = \tau_i$ . By SMLRP we have:

$$P(Y = 1)\frac{f_{\Psi(X)|Y=1}(\tilde{\tau}_1)}{f_{\Psi(X)}(\tilde{\tau}_1)} = \left(1 + \frac{1-p}{p}\frac{1}{\Lambda_\Psi(\tilde{\tau}_1)}\right)^{-1} <$$

$$< \left(1 + \frac{1-p}{p}\frac{1}{\Lambda_\Psi(\tau_1)}\right)^{-1} = P(Y = 1)\frac{f_{\Psi(X)|Y=1}(\tau_1)}{f_{\Psi(X)}(\tau_1)} < r_2$$

Therefore (11) is maintained (the other inequalities are not affected). On the other hand by theorem A.1 we have strict monotonicity of $R$, meaning:

$$R_1(\Psi, \tilde{\tau}) = P(Y = 1 \mid \Psi(X) < \tilde{\tau}_1) < P(Y = 1 \mid \Psi(X) < \tau_1) = R_1(\Psi, \tau) = r_1$$

and therefore $\mathrm{IRD}_r(\Psi, \tilde{\tau}) > 0$. The conclusion is that even under SMLRP we can use (11),(12) only as necessary conditions for the feasibility of a given solution and that the test of feasibility must be performed using (1) and (2) directly.

In order to satisfy the necessary conditions for IRD=0 in the absence of SMLRP we can generally require $r_i > \inf\limits_{\{\tau_i : \tau_i > \tau_{i-1}\}} R_i(\Psi, \tau)$ (we require strong inequalities to avoid degenerate zero-length intervals), however for such cases the existence of a closed-form expression would depend on the exact distribution of $X|Y = k$ ($k = 0, 1$). We leave the exact formulation of non-SMLRP lower bounds outside the scope of this paper.

The final conclusion is that given two sets of risk categories $r_1, r_2$ and a couplet $(\Psi, \tau_1)$ which satisfies $\mathrm{IRD}_{r_1}(\Psi, \tau_1) = 0$, we may not be able to find a set of breakpoints $\tau_2$ which satisfies $\mathrm{IRD}_{r_2}(\Psi, \tau_2) = 0$ (using the same model $\Psi$). Specifically we can now claim that optimal models of existing classification methods (typically optimized for $r = (0, 1)$) are not necessarily feasible for any choice of $r$.

The existence of lower bounds on the IRD is perhaps the most counter-intuitive result of this paper. The reason why these limitations have not been addressed before has to do with the fact that most classification methods use a single breakpoint to distinguish between the two groups ($\tau \in \mathbb{R}$) and the issue of degenerate solutions or non-feasibility of $\Psi$ is avoided altogether. Although the fulfilment of (11),(12) does not ensure the feasibility of a given solution, these inequalities are instrumental in demonstrating why the solutions from existing methods may not be feasible for a different choice of $r$, and provide an elegant method to disqualify such solutions. Once we define our objective as the distribution pre-set risk levels over multiple adjacent intervals we must recognize the existence of possible limitations on IRD for existing methods and as a result define new conditions for optimality.

# 4 Ordinal Risk-Group Classification

Although the definition of IRD naturally suggests itself as a new criterion for optimality (look for a couplet $(\Psi, \tau)$ such that $\mathrm{IRD}_r(\Psi, \tau(\Psi)) = 0$), there are two problems with using IRD as a single optimality criterion. First, since our problem is a classification problem we must consider some sense of the quality of separation between the two classes in order to avoid degenerate solutions. This principle is not straight forwardly reflected by the definition of IRD (2). Second, our definition of IRD and the resulting necessary inequalities (11) do not ensure existence or uniqueness of an optimal solution.

Our practical solution to these problems is to define IRD as a feasibility criterion and use it as a constraint in an existing optimization problem. Since we are still in the domain of classification problems it would be reasonable to preserve some basic concepts, particularly the definition of optimality: We seek a model that on the one hand maximizes our ability to discriminate between the two classes, but on the other hand distributes risk correctly, meaning that it belongs to the set of feasible solutions:

$$C_r(0) = \{(\Psi, \tau) : \mathrm{IRD}_r(\Psi, \tau) = 0\} \tag{13}$$

In the event that $C$ is an empty set we would have to reconsider our pre-set $r$ or change our method of constructing $\Psi$.

Any classification method we might consider for IRD "augmentation" must satisfy several criteria. First, it must provide a continuous output $\Psi(X)$, ensuring that we have an appropriate output that can be partitioned into intervals (using $\tau$). This requirement automatically excludes classification methods that do not combine the vector of explanatory variables $X$ into a single real-valued score $\Psi(X)$ before making a prediction of risk or class membership (classification trees are an example of such excluded methods). Furthermore, we would like to maintain the notion that observations with higher scores have a higher conditional risk, and therefore require that the output $\Psi(X)$ is strongly correlated with the conditional risk function $P(Y = 1 \mid \Psi(X) = x)$. Methods such as Fisher's LDA [8] or SVM for two classes do provide a continuous scale and a single breakpoint to predict class membership, however these scales are not necessarily correlated with the conditional risk and only ensure that a majority of the observations from the special class are on one side of the breakpoint. We therefore decided to focus our discussion on risk estimation methods that provide a direct estimation of the risk function:

$$\Psi : \mathbb{R}^P \longrightarrow [0, 1], \quad \Psi(X) = P(Y = 1 \mid X) \tag{14}$$

Finally, in order to simplify our construction we assume SMLRP (see appendix A). As we have seen before, this assumption ensures that we have a simple way to calculate the lower bounds on IRD, and also ensures that for a given model $\Psi$, if exists $\tau(\Psi)$ such that $\mathrm{IRD}_r(\Psi, \tau(\Psi)) = 0$ then it is unique (see lemma B.1). These properties enable us to simplify our parameter space by optimizing over $\Psi$ alone, and provide a simple way to test for the existence of necessary conditions for $\mathrm{IRD}_r(\Psi) = \mathrm{IRD}_r(\Psi, \tau(\Psi)) = 0$ and optimize under the constraint $C_r(\Psi) = \{\Psi : \mathrm{IRD}_r(\Psi) = 0\}$.

## 4.1 Penalized Optimization

While the idea of fitting an optimal risk predictor that maximizes class discrimination is a well defined concept, the requirement of $\text{IRD}_r(\Psi) = 0$ may lead to degenerate solutions of $\tau(\Psi)$ for certain values of $r$. We demonstrate this problem for a simple case of homoscedastic one-dimensional Gaussian logistic regression: Let $X \mid Y = 1 \sim N(\mu, \sigma)$, $X \mid Y = 0 \sim N(-\mu, \sigma)$, $P(Y = 1) = \frac{1}{2}$ and the model $\Psi(\beta, x)$ is the one-dimensional logistic function with the parameter $\beta$, meaning $\Psi : \mathbb{R} \times \mathbb{R} \longrightarrow [0, 1]$, $\Psi(\beta, x) = \frac{\exp(\beta x)}{1 + \exp(\beta x)}$. We set $r = (0, 0.5, 1)$.

Denoting $\tau(\Psi, \beta, t) = (\Psi(\beta, -t), \Psi(\beta, t))$, we use symmetry of the conditional distributions around $x = 0$ and the strict monotonicity of $\Psi(\beta, x)$ in $x$ and $\beta$ to show that for any choice of $\beta, t \in \mathbb{R}$ we can minimize error for $i = 2$:

$$R_2(\Psi(\beta, X), \tau(\Psi, \beta, t)) = P(Y = 1 \mid \Psi(\beta, X) \in (\Psi(\beta, -t), \Psi(\beta, t)])$$
$$= P(Y = 1 \mid \beta X \in (-\beta t, \beta t]) = P(Y = 1 \mid X \in (-t, t]) = 0.5 = r_2$$

Similar considerations ensure that the risk prediction errors are equal on both sides:

$$R_1(\Psi(\beta, X), \tau(\Psi, \beta, t)) = P(Y = 1 \mid X < -t)$$
$$= 1 - P(Y = 1 \mid X > -t) = 1 - R_3(\Psi(\beta, X), \tau(\Psi, \beta, t))$$

Using Bayes theorem we have:

$$P(Y = 1 \mid X < -t) = \frac{P(Y = 1)P(X < -t \mid Y = 1)}{P(Y = 1)P(X < -t \mid Y = 1) + P(Y = 0)P(X < -t \mid Y = 0)}$$
$$= \frac{P(Y = 1)\Phi(-t - \mu)}{P(Y = 1)\Phi(-t - \mu) + P(Y = 0)\Phi(-t + \mu)}$$
$$= \left(1 + \frac{P(Y = 0)}{P(Y = 1)} \frac{\Phi(-t + \mu)}{\Phi(-t - \mu)}\right)^{-1}$$

where $\Phi$ is the CDF of the standard normal distribution. Using the known inequality:

$$\frac{\phi(x)}{x + 1/x} < \Phi(-x) < \frac{\phi(x)}{x} \quad \forall x > 0, \tag{15}$$

where $\phi$ is the PDF of the standard normal distribution, we show an upper bound:

$$\frac{\Phi(-t + \mu)}{\Phi(-t - \mu)} > \frac{(t + \mu) + \frac{1}{t+\mu}}{t - \mu} \frac{\phi(t - \mu)}{\phi(t + \mu)} = \frac{(t + \mu) + \frac{1}{t+\mu}}{t - \mu} e^{2\mu t} \quad \forall t > \mu$$

Therefore $\lim_{t \to \infty} P(Y = 1 \mid X < -t) = 0$ and similarly $\lim_{t \to \infty} P(Y = 1 \mid X > t) = 1$. For any arbitrarily small $\varepsilon > 0$ we can find a sufficiently large $t$ such that

$$R_1(\Psi(\beta, X), \tau(\Psi, \beta, t)) = P(Y = 1 \mid X < -t) = \leq \varepsilon/2$$

making the total IRD:

$$\text{IRD}_r(\Psi(\beta, X), \tau(\Psi, \beta, t)) = \sqrt{\sum_{i=1}^{3} (R_i(\Psi(\beta, X), (\Psi(\beta, -t), \Psi(\beta, t))) - r_i)^2} \leq \varepsilon$$

10

As a result, for any given $\beta$ the only solution that satisfies IRD=0 is degenerate:

$$\lim_{t \to \infty} \text{IRD}_r(\Psi(\beta, X), \tau(\Psi, \beta, t)) = 0$$

There are several alternatives for dealing with this problem. First, we may decide that methodologically we do not allow setting $r_1 = 0$ or $r_T = 1$. This will ensure that the values of $\tau$ are finite but might still lead to very large or very small intervals, depending on the parameters of the model. Alternatively, if our risk estimation method uses optimization to fit the optimal model (for example maximizing the likelihood function in the case of parametric methods) then we can introduce a penalty function $\text{Pen} : \mathbb{R}^{T-1} \to \mathbb{R}$, which will enable us to balance the properties of $\tau$ (minimal or maximal distance between breakpoints) with the discrimination properties of $\Psi$. This means that instead of maximizing or minimizing a target function $f(\Psi \mid X, Y)$ we maximize/minimize $f(\Psi \mid X, Y) + \gamma \text{Pen}(\tau)$ under an IRD constraint, where $\gamma$ is a tuning parameter that represents the degree of aversion to degenerate solutions.

In cases where the degenerate solutions are encountered we would opt for the use of a penalty function. This reflects our understanding that the requirement of "evenly spread" breakpoints is relatively subjective and should allow for some discretion as to the balance between the ability of the model to separate classes and the resulting interval lengths. By choosing an appropriate penalty function and an aversion parameter $\gamma$ we enable better fitting of the model according to the circumstances at hand, while introducing a relatively small number of additional parameters. On the other hand, since IRD represents an absolute measure of the model's quality, we believe it must be tightly controlled as the constraint $\text{IRD}_r(\Psi, \tau) = 0$ on any model we might consider. We address the details of constructing this constraint for parametric models in the following section.

## 4.2  Estimation of Interval Risk Deviation

So far, we have defined interval risk deviation (IRD) as a property of a score model $\Psi$ and the joint distribution of $(X, Y)$. In order to implement the concept of IRD in a real-life scenario we must describe a way to estimate $\text{IRD}_r(\Psi)$ based on a sample of $N$ i.i.d observations from a known $P$-dimensional multivariate distribution $\mathcal{F}(\theta)$ in the form of a $N \times P$ matrix $\mathbf{X}$ and a vector $y \in \{0, 1\}^N$ representing known class memberships (depending on the design of the experiment $y$ may or may not be a random sample). Focusing on parametric methods, we assume that $P(Y = 1 \mid X) = \Psi(\beta, X)$ where $\Psi : \mathbb{R}^M \times \mathbb{R}^P \to [0, 1]$ is a known function and $\beta \in \mathbb{R}^M$ is the set of parameters controlling the shape of the function (e.g. generalized linear models [20] where $M = P$, $g$ is a known, strictly monotone and bijective link function and $\Psi(\beta, x) = g(\beta^T x)$). Having previously assumed SMLRP and a closed-from $\Psi$, we can simplify our notation by denoting $\tau(\beta) = \tau(\Psi(\beta, X))$, IRD for a given $\beta$ as $\text{IRD}_r(\beta) = \text{IRD}_r(\Psi(\beta, X), \tau(\Psi(\beta, X)))$ and the constraint set $C_r(\beta) = \{\beta \in \mathbb{R}^M : \text{IRD}_r(\beta) = 0\}$

Many parametric classification methods solve the problem of estimating $\beta$ from a given sample $(\mathbf{X}, y)$ by using the maximum likelihood (ML) method. We denote

$\mathbf{X}_{j,\cdot}$ the $j$'th row of the matrix $\mathbf{X}$, making our model-predicted probability for the $j$'th observation $\Psi(\beta, \mathbf{X}_{j,\cdot}) = P(Y = 1 \mid \mathbf{X}_{j,\cdot})$. Assuming random sampling, there are two equivalent formulations of the "complete" likelihood function:

$$L(\theta_0, \theta_1, p \mid \mathbf{X}, y) = \prod_{j=1}^{N} f_{X,Y}(\mathbf{X}_{j,\cdot}, y_j) = \prod_{j=1}^{N} f_{\theta_{y_j}}(\mathbf{X}_{j,\cdot}) P(Y_j = y_j) \tag{16}$$

$$L(\beta, \theta \mid \mathbf{X}, y) = \prod_{j=1}^{N} P(Y = y_j \mid \mathbf{X}_{j,\cdot}) f_\theta(\mathbf{X}_{j,\cdot}) \tag{17}$$

where $f_\theta$ is the density function of the distribution $\mathcal{F}(\theta)$ of $X$ and $f_{\theta_{y_j}}$ is the density function of the conditional distribution $\mathcal{F}_{y_j}(\theta_{y_j})$ of $X \mid Y = y_j$. When using (16) we must make additional assumptions about the conditional distribution of $X \mid Y = k$ and a random sampling process, and as a result our estimator $\hat{\beta}$ becomes a function of the estimators $\hat{\theta}_0, \hat{\theta}_1, \hat{p}$. On the other hand using (17) can significantly simplify the optimization process. Since our parameter of interest $\beta$ is isolated in the term $P(Y = 1 \mid \mathbf{X}_{j,\cdot}) = \Psi(\beta, \mathbf{X}_{j,\cdot})$ and does not effect the term $f_\theta(\mathbf{X}_{j,\cdot})$, we can directly maximize the the partial likelihood function $L_\Psi$ over the values of $\beta$:

$$L_\Psi(\beta \mid \mathbf{X}, y) = \prod_{j=1}^{N} P(Y = y_j \mid \mathbf{X}_{j,\cdot}) = \prod_{j=1}^{N} \Psi(\beta, \mathbf{X}_{j,\cdot})^{y_j} (1 - \Psi(\beta, \mathbf{X}_{j,\cdot}))^{1-y_j} \tag{18}$$

making the corresponding maximum likelihood optimization problem:

$$\hat{\beta}_{ML} = \underset{\beta \in \mathbb{R}^M}{\operatorname{argmax}} \, L_\Psi(\beta \mid \mathbf{X}, y) \tag{19}$$

One of the primary advantages of using the second approach for maximum likelihood estimation is that it circumvents the need to estimate the parameters of the conditional distributions, thus making $\beta$ the only estimated parameter. This construction also enables a relatively simple extension of the maximum likelihood framework to semi-parametric models or non-random sampling (see [23] for such an extension of logistic regression).

Incorporating the IRD constraint into the maximum likelihood framework means that for any given $\beta$ we must be able to estimate the conditional probability over intervals $R_i(\beta, \tau) = P(Y = 1 \mid \Psi(\beta, X) \in (\tau_{i-1}, \tau_i])$ for an arbitrary $\tau$ from the sample $(\mathbf{X}, y)$, which we can then use to calculate the estimate for the unique optimal $\tau(\beta)$ (which is a function of both $\beta$ and the distribution of $X$). Alas, it is usually difficult to derive $R_i(\beta, \tau)$ directly from the point-wise conditional probability $P(Y = 1 \mid X)$. In order to facilitate the estimation of IRD for such cases we use an approach similar to (16) but in a different context. As we shall see this will enable us to provide parametric estimates of IRD while retaining the convenient structure of our target function $L_\Psi$.

We assume that $\Psi(\beta, X) \mid Y = k$ has a known distribution which we denote $F(\eta_k(\beta))$, and that $F$ is continuous. We can now use Bayes's theorem to represent $R(\beta, \tau)$ as:

$$R(\beta, \tau) = R(\Psi(\beta, X), \tau) = \left(1 + \frac{1-p}{p} \frac{1}{\nu(\beta, \tau)}\right)^{-1} \tag{20}$$

where:

$$\nu_i(\beta, \tau) = \frac{P(\Psi(\beta, X) \in (\tau_{i-1}, \tau_i] \mid Y = 1)}{P(\Psi(\beta, X) \in (\tau_{i-1}, \tau_i] \mid Y = 0)} = \frac{F_{\eta_1(\beta)}(\tau_i) - F_{\eta_1(\beta)}(\tau_{i-1})}{F_{\eta_0(\beta)}(\tau_i) - F_{\eta_0(\beta)}(\tau_{i-1})} \tag{21}$$

Therefore by estimating the parameters $p$, $\eta_0(\beta)$, $\eta_1(\beta)$ we can calculate the estimators $\hat{\nu}(\beta, \tau)$ and $\hat{R}(\beta, \tau)$ and use them to find the unique $\hat{\tau}(\beta)$ which solves:

$$\hat{\mathrm{IRD}}_r(\beta) = \mathrm{IRD}_r(\beta, \hat{\tau}(\beta)) = 0 \tag{22}$$

The complete likelihood function under these assumptions is:

$$L(\eta_0(\beta), \eta_1(\beta), p \mid \mathbf{X}, y, \beta) = \prod_{j=1}^{N} f_{\eta_{y_j}(\beta)}(\Psi(\beta, \mathbf{X}_{j,\cdot})) P(Y_j = y_j) \tag{23}$$

where $f_{\eta_{y_j}(\beta)}$ is the density function of $\Psi(X) \mid Y = y_j \sim F(\eta_{y_j}(\beta))$. Since the estimation of $p$, $\eta_0(\beta)$, $\eta_1(\beta)$ is performed for each $\beta$ independently and only in order to verify the compliance with the constraint $\hat{\mathrm{IRD}}_r(\beta) = 0$, the estimation of these parameters does not effect the value of the target function $L_\Psi(\beta \mid \mathbf{X}, y)$. We use this property together with the separability of the likelihood function to estimate $p$ and $\eta_0(\beta), \eta_1(\beta)$ separately for each beta. For the case of random sampling, the parameter $p$ is relatively easy to estimate independently of $\beta$ as:

$$\hat{p} = \hat{P}(Y = 1) = \frac{\sum_{j=1}^{N} y_i}{N} \tag{24}$$

although for other cases, such as case-control studies, we might need additional information to correct for biased sampling. For the estimation of $\eta_0(\beta), \eta_1(\beta)$ we use our assumptions about $F$ to build an ancillary optimization problem for each $\beta$ individually, and use the estimates provided by the ancillary problem (conditional on the value of $\beta$) to estimate $\hat{\mathrm{IRD}}_r(\beta, \tau)$ (again the target function $L_\Psi$ remain unchanged). The likelihood function for the ancillary problem can be rewritten as:

$$L_F(\eta_0(\beta), \eta_1(\beta) \mid \mathbf{X}, y, \beta) = \prod_{\{j\,:\,y_j=0\}} f_{\eta_0(\beta)}(\Psi(\beta, \mathbf{X}_{j,\cdot})) \prod_{\{j\,:\,y_j=1\}} f_{\eta_1(\beta)}(\Psi(\beta, \mathbf{X}_{j,\cdot})) \tag{25}$$

where $f_{\eta_k(\beta)}$ is the density function of the distribution $\Psi(\beta, X) \mid Y = k \sim F(\eta_k(\beta))$, and the ancillary maximum likelihood estimation problem is:

$$(\hat{\eta}_0(\beta), \hat{\eta}_1(\beta)) = \operatorname*{argmax}_{\eta_0(\beta), \eta_1(\beta)} L_F(\eta_0(\beta), \eta_1(\beta) \mid \mathbf{X}, y, \beta) \tag{26}$$

The construction of multiple ancillary ML estimation problems can be computationally demanding, but luckily for many known distributions the formula for the ML estimator $\hat{\eta}_k(\beta)$ is known, and in fact it is relatively simple to calculate it directly from the known ML estimators of the conditional distribution $X \mid Y = k \sim \mathcal{F}(\theta_k)$ as $\hat{\eta}_k(\beta) = \eta(\beta, \hat{\theta}_k)$. This fact significantly reduces the complexity of estimating the IRD constraint. For a concrete example of such a case using Gaussian logistic regression (which can be easily extended to other GLM instances) see section 5.

Alternatively it would be possible to use non-parametric estimators or approximation. While it is possible to attempt to directly approximate $Q_i(\beta) = P(Y = 1 \mid \Psi(\beta, X) \leq \tau_i(\beta))$, $P(\Psi(\beta, X) \leq \tau_{i-1}(\beta))$, $P(\Psi(\beta, X) \leq \tau_i(\beta))$ and utilize (5) for a direct estimation of $R(\beta)$ (in this case the equality is useful since $\beta$ is fixed), it would often be more convenient to approximate $p$ and $F_{\eta_k(\beta)}$ at $\{\tau_i(\beta)\}_{i=1}^{T-1}$ (a total of $2T - 1$ approximations per $\beta$) and use (20) to calculate $\hat{\nu}(\beta)$ and the resulting IRD estimate. The primary advantage of this approach is that it requires no additional assumptions about the distribution of $(X, Y)$ and can therefore be easily extended to other non-ML estimation methods of $\beta$. On the other hand, by using non-parametric methods we pay a price both in the quality of our estimates (we ignore information about the distribution of $X$) and in the computational complexity of our estimation scheme.

Finally, regardless of our approach to estimation, we recognize the fact that under realistic scenarios we will have to use numeric methods for the calculation of $\hat{\tau}(\beta)$ and for the estimation of the required parameters. We therefore set a low threshold $\varepsilon$ and accept $\beta$ as feasible if our estimated IRD satisfies:

$$\hat{\mathrm{IRD}}_r(\beta) = \mathrm{IRD}_r(\beta, \hat{\tau}(\beta)) = \|\hat{R}(\beta, \hat{\tau}(\beta)) - r\| < \varepsilon \qquad (27)$$

making our feasibility set $\hat{C}_r(\varepsilon) = \{\beta \in \mathbb{R}^M : \hat{\mathrm{IRD}}_r(\beta) < \varepsilon\}$ and the penalized Ordinal Risk-Group (ORG) optimization problem:

$$\hat{\beta}_{ORG} = \underset{\beta \in \hat{C}_r(\varepsilon)}{\mathrm{argmax}} \, L_{\Psi}(\beta \mid \mathbf{X}, y) + Pen(\hat{\tau}(\beta)) \qquad (28)$$

If for our choice of $r$ we have $\hat{\beta}_{LR} \in \hat{C}_r(\varepsilon)$ and the distances between the set of breakpoints $\hat{\tau}(\hat{\beta}_{LR})$ are non-degenerate, then the global (unconstrained) optimality of $\hat{\beta}_{LR}$ ensures that it is also the optimal solution of the constrained ordinal problem. It may also serve as the optimal solution for the constrained and penalized ordinal problem, but that will depend on the selection of the aversion parameter. On the other hand, as we've seen in section 3.2, once we have more than a single breakpoint we introduce limits of feasibility into the maximization problem and may discover that the solution to (36) is no longer feasible ($\hat{\mathrm{IRD}}_r(\hat{\beta}_{LR}) \geq \varepsilon$). For such cases we must define a new constrained optimization problem and look for a new optimal solution. We discuss an example for such a case in the following section.

There are two issues we leave outside the scope of this paper. First, although the consistency of constrained ML estimation has been explored in various contexts (for example for mixture models [11]), the consistency of the ML estimators under the specific constraint of $\hat{\mathrm{IRD}} = 0$ requires verification. Similarly, although in our description of the problem $\tau$ is a function of $\beta$ and the parameters of $\mathcal{F}$ (a result of the uniqueness demonstrated in appendix B), it remains to be verified whether the consistency of the estimator $\hat{\beta}$ ensures the consistency of $\hat{\tau}(\hat{\beta})$. We expect the fact that for many cases $\hat{\tau}(\beta)$ has no analytical solution to further complicate this problem.

Second, in order to measure our estimation errors we require a method for building right-sided confidence intervals for IRD based on the distribution of $\hat{\mathrm{IRD}}$ for a given $\beta$. Since $\nu(\beta, \tau(\beta))$ is a ratio of CDF differences, the process of deriving the

distribution of $\hat{\nu}(\beta, \hat{\tau}(\beta))$ from the distribution of $\hat{\eta}_k(\beta)$ would require several steps of approximation, primarily since $\hat{\tau}(\beta)$ the result of numeric estimation (even if the distribution of $\hat{\eta}_k(\beta)$ is known). Similar problems apply for non-parametric estimators, although we can see two possible approached for a solution. The first approach would be to use an equivalent definition of IRD as $\text{IRD}_r(\Psi, \tau) = \max_i |R_i(\Psi, \tau) - r_i|$ and try to prove a Glivenko-Cantelli [28] type theorem for conditional distributions which would describe the necessary conditions ensuring that:

$$\sup_{x_2 > x_1} |\hat{P}_N(Y = 1 \mid X \in (x_1, x_2]) - P(Y = 1 \mid X \in (x_1, x_2])| \xrightarrow[N \to \infty]{} 0 \qquad (29)$$

where $\hat{P}_N$ is the empirical conditional probability estimator:

$$\hat{P}_N(Y = 1 \mid X \in (x_1, x_2]) = \frac{|\{j \ : \ \Psi(\mathbf{X}_j) \in (x_1, x_2] \ \wedge \ y_j = 1\}|}{|\{j \ : \ \Psi(\mathbf{X}_j) \in (x_1, x_2]\}|} \qquad (30)$$

Building on this result we can attempt to derive the asymptotic distribution of (29) and try to construct test that will be the conditional equivalent of the Kolmogorov-Smirnov test [28]. The second approach, which is less elegant but more plausible, would be to combine (21) with the well known asymptotic behaviour of the empirical conditional distribution function:

$$\hat{F}_{\eta_k(\beta)}^{(N_k)}(t) = \frac{|\{j \mid \Psi(\beta, \mathbf{X}_{j,\cdot}) \leq t, y_j = k\}|}{N_k} \quad k = 0, 1 \qquad (31)$$

where $N_k = |\{j : y_j = k\}|$. By the central limit theorem [28], this estimator weakly converges to $F_{\eta_k(\beta)}$ pointwise:

$$\sqrt{N_k} \left( \hat{F}_{\eta_k}^{(N_k)}(t) - F_{\eta_k}(t) \right) \xrightarrow[N_k \to \infty]{} N\left(0, F_{\eta_k}(t)\left(1 - F_{\eta_k}(t)\right)\right) \qquad (32)$$

However the fact that $\hat{\tau}(\beta)$ changes as a function of the sample (and is not a fixed quantile $t$) means that the points where $\hat{F}_{\eta_k}^{(N_k)}$ is estimated change as a function of the sample ($N_k$), therefore the nature of the convergence and the resulting asymptotic distribution depend on the convergence of $\hat{\tau}(\beta) \to \tau(\beta)$. We would therefore need to find sufficient conditions for:

$$\tau(\hat{\beta})^{(N_k)} \to \tau(\beta) \Rightarrow$$
$$\sqrt{N_k} \left( \hat{F}_{\eta_k}^{(N_k)}(\tau(\hat{\beta})^{(N_k)}) - F_{\eta_k}(\tau(\hat{\beta})^{(N_k)}) \right) \xrightarrow[N_k \to \infty]{} N\left(0, F_{\eta_k}(\tau(\beta))\left(1 - F_{\eta_k}(\tau(\beta))\right)\right)$$
$$(33)$$

The next step would be to use the strong uniform convergence of $\hat{F}_{\eta_k}^{(N_k)} \to \hat{F}_{\eta_k}$ (as provided by the Glivenko-Cantelli theorem) to approximate $\hat{F}_{\eta_k}^{(N_k)}(\hat{\tau}_i(\beta))$ as normally distributed $\mu = \hat{F}_{\eta_k}^{(N_k)}(\hat{\tau}_i(\beta))$ and $\sigma^2 = \hat{F}_{\eta_k}^{(N_k)}(\hat{\tau}_i(\beta))\left(1 - \hat{F}_{\eta_k}^{(N_k)}(\hat{\tau}_i(\beta))\right)$. Combined with Donsker's theorem [6], we should be able to estimate the asymptotic distribution of the difference of two points of $\hat{F}_{\eta_k}^{(N_k)}$ (which make both the numerator and the denominator of $\hat{n}u_i$) as normallt distributed with mean $\mu = \hat{F}_{\eta_k}(\hat{\tau}_i(\beta)) - \hat{F}_{\eta_k}(\hat{\tau}_{i-1}(\beta))$ and variance $\sigma^2 = \left(\hat{F}_{\eta_k}(\hat{\tau}_i(\beta)) - \hat{F}_{\eta_k}(\hat{\tau}_{i-1}(\beta))\right)\left(1 - \hat{F}_{\eta_k}(\hat{\tau}_i(\beta)) - \hat{F}_{\eta_k}(\hat{\tau}_{i-1}(\beta))\right)$. The

final step would be to use work by Hinkley [12], which describes the distribution of a ratio of two non-correlated normal random variables, to approximate the distribution of $\hat{\nu}_i(\beta, \hat{\tau}(\beta))$ which we can use to estimate $P(\hat{\text{IRD}}_r(\beta) > \varepsilon)$. We note that the construction of such test would also mean that we can change the definition of our feasibility set to $\hat{C}_r(\varepsilon, \alpha) = \{\beta \in \mathbb{R}^M : P(\hat{\text{IRD}}_r(\beta) < \varepsilon) > 1 - \alpha\}$. We leave the details and proof of these ideas for future papers.

# 5 Case study: Gaussian Logistic Regression

Logistic regression is one of the most studied classification methods in the scientific literature and has been widely applied in statistics, scientific research and industry. The name "logistic" for the function $f(x) = \frac{e^x}{1+e^x}$ was originally coined by Verhulst as early as the 19'th century, but it was Cox [4] who used it first in the context of binary data analysis. The concept of multinomial logistic regression was first suggested by Cox (1966) [3] and developed independently by Theil (1969) [26]. The link to ordered choice models (ordered logistic regression) was made by McFadden in his paper from 1973 [19]. Cramer (2002) [5] has a complete historical review.

Logistic regression belongs to the group of classification methods that estimate class membership probability rather than predict class membership. It is a special instance of a larger group of parametric models called generalized linear models (GLM [20]), which extends linear models by allowing the addition of a predefined link function $g$ that connects the linear model $\beta^T X$ ($\beta \in \mathbb{R}^P$) to the response variable $Y$ (meaning that $g^{-1}(Y) = \beta^T X$) and assuming that the distribution of $X$ is from an exponential family. In the case of logistic regression the link function is assumed to be the logistic function $g(\beta^T x) = \frac{e^{\beta^T x}}{1+e^{\beta^T x}}$, making the inverse function $g^{-1}(p) = \text{logit}(p) = log\left(\frac{p}{1-p}\right)$. The probability of belonging to the "special class" (in the case of 2 classes) conditioned on the r.v $X$ is assumed to be:

$$P_\beta(Y = 1 \mid X) = \frac{e^{\beta^T X}}{1 + e^{\beta^T X}} \tag{34}$$

where the assumptions on $(X, Y)$ can be modified to match a wide variety of cases, for example to a non-random sampling scheme like case-control studies or semi-parametric models [23].

For the purpose of demonstrating our method we assume that $\mathbf{X}$ is a $N \times P$ matrix representing $N$ i.i.d random samples from a $P$-dimensional multivariate normal distribution. Under this assumption the vector $y \in \{0, 1\}^N$ of class memberships represents the result of $N$ independent Bernoulli random variables $\{Y_j\}_{j=1}^T$. Even under these assumptions, the ML problem does not have an analytical least squares solution, and is usually solved using numerical maximum likelihood algorithms. The log-likelihood function is:

$$l_{LR}(\beta \mid \mathbf{X}, y) = \log(L_{LR}(\beta \mid \mathbf{X}, y)) = \sum_{j=1}^N y_j \beta^T \mathbf{X}_{j,\cdot} - \sum_{j=1}^N \log\left(1 + e^{\beta^T \mathbf{X}_{j,\cdot}}\right) \tag{35}$$

16

and the logistic regression (LR) maximum likelihood optimization problem is:

$$\hat{\beta}_{LR} = \operatorname*{argmax}_{\beta \in \mathbb{R}^P} l_{LR}(\beta \mid \mathbf{X}, y) \tag{36}$$

As we have noted in section 4.2, the parametric estimation of IRD requires several additional assumptions. We assume that the conditional distributions $X \mid Y = k$ ($k = 0, 1$) are also multi-variate normal, and in order to achieve SMLRP we assume equal conditional covariance. Using terms defined to construct (28) our assumptions on Gaussian logistic regression translate into the following:

$$\Psi(\beta, x) = \frac{\exp(\beta^T x)}{1 + \exp(\beta^T x)}, \quad X \mid Y = k \sim MVN(\underline{\mu}_k, \Sigma) \tag{37}$$

As a result:

$$\operatorname{logit}(\Psi(\beta, X)) \mid Y = k \sim N(\mu_k(\beta) = \beta^T \underline{\mu}_k, \sigma^2(\beta) = \beta^T \Sigma \beta) \tag{38}$$

Conveniently, the ML estimators follow a similar pattern. For the construction of the known ML estimators of the distribution of $X$ we denote $\mathbf{X}^{(k)}$ as the matrix composed of all the lines of $\mathbf{X}$ for which $y_j = k$, $N_1 = \sum_{j=1}^N y_j$, $N_0 = N - \sum_{j=1}^N y_j$ and $\overline{\mathbf{X}}_m^{(k)} = \frac{\sum_{j=1}^{N_k} \mathbf{X}_{j,m}^{(k)}}{N_k}$ as the average of the $m$'th column of $\mathbf{X}^{(k)}$ ($m = 1, \ldots, P$). The ML estimators are:

$$\hat{\underline{\mu}}_k = \overline{\mathbf{X}}^{(k)} = (\overline{\mathbf{X}}_1^{(k)}, \ldots, \overline{\mathbf{X}}_P^{(k)}) \quad (k = 0, 1) \tag{39}$$

and having assumed equal covariance we use the pooled covariance matrix estimator:

$$\hat{\Sigma}_k = \frac{1}{N_k} \sum_{j=1}^{N_k} (\mathbf{X}_{j,\cdot}^{(k)} - \overline{\mathbf{X}}^{(k)})^T (\mathbf{X}_{j,\cdot}^{(k)} - \overline{\mathbf{X}}^{(k)})$$

$$\hat{\Sigma} = \frac{1}{N}(N_0 \hat{\Sigma}_0 + N_1 \hat{\Sigma}_1) \tag{40}$$

The assumption of multivariate-normal conditional distributions enables us to avoid the construction of an ancillary ML problem for each $\beta$ by using the known relationship between the ML estimators $(\hat{\underline{\mu}}_0, \hat{\underline{\mu}}_1, \hat{\Sigma})$ and the $\beta$-transformed ML estimators $(\hat{\mu}_k(\beta), \hat{\sigma}(\beta))$:

$$\hat{\mu}_k(\beta) = \beta^T \hat{\underline{\mu}}_k, \quad \hat{\sigma}(\beta) = \sqrt{\beta^T \hat{\Sigma} \beta} \tag{41}$$

The strict monotonicity of $logit(p) = log\left(\frac{p}{1-p}\right) = \Psi_\beta^{-1}(p)$ in $p$ means that we can use the equality:

$$\begin{aligned} &P(\Psi(\beta, X) \in (\tau_{i-1}(\beta), \tau_i(\beta)] \mid Y = k) \\ &= P(\beta^T X \in (\operatorname{logit}(\tau_{i-1}(\beta)), \operatorname{logit}(\tau_i(\beta))] \mid Y = k) \\ &= \Phi\left(\frac{\operatorname{logit}(\tau_i(\beta)) - \mu_k(\beta)}{\sigma(\beta)}\right) - \Phi\left(\frac{\operatorname{logit}(\tau_{i-1}(\beta)) - \mu_k(\beta)}{\sigma(\beta)}\right) \end{aligned} \tag{42}$$

to estimate $\nu_i(\beta)$ as:

$$\hat{\nu}_i(\beta) = \frac{\Phi\left(\frac{\text{logit}(\hat{\tau}_i(\beta)) - \hat{\mu}_1(\beta)}{\hat{\sigma}(\beta)}\right) - \Phi\left(\frac{\text{logit}(\hat{\tau}_{i-1}(\beta)) - \hat{\mu}_1(\beta)}{\hat{\sigma}(\beta)}\right)}{\Phi\left(\frac{\text{logit}(\hat{\tau}_i(\beta)) - \hat{\mu}_0(\beta)}{\hat{\sigma}(\beta)}\right) - \Phi\left(\frac{\text{logit}(\hat{\tau}_{i-1}(\beta)) - \hat{\mu}_0(\beta)}{\hat{\sigma}(\beta)}\right)} \tag{43}$$

Finally, we use the assumption of random sampling to estimate $\hat{p} = \frac{1}{N}\sum_{j=1}^{N} y_j$ and utilize (20) to construct our parametric estimation of IRD for logistic regression.

## 5.1 Example: The Wisconsin Diagnostic Breast Cancer (WDBC) Dataset

In this section we bring an example of the sub-optimality of using one of the most commonly used classification methods - Logistic Regression (LR) - to solve a relatively simple ordinal problem. We then provide the Ordinal Risk-Group version of Logistic Regression (ORG-LR) solution to the problem and compare our results.

The dataset we used for this example is the extensively used Wisconsin Diagnostic Breast Cancer [25] dataset from the UCI Machine Learning Repository [9], which contains the analysis of cell nuclei from 556 patients using digitized images of fine needle aspirate (FNA) of extracted breast masses. Since we assumed a continuous $\mathcal{F}$ and equal variance we selected the following features for the construction of our models: texture, log area, smoothness, log compactness, log concave points, log symmetry and an intercept variable. The final result from the diagnosis ("Malignant" $N = 212$/ "Non-Malignant" $N = 344$) was used as the dependent variable for the logistic regression analysis and ordinal risk group analysis. The code for this example was written in the R programming language [24] using internal optimization algorithms and the Augmented Lagrangian Adaptive Barrier Minimization Algorithm (the alabama library) with the constraint $\hat{\text{IRD}}_r(\beta) < \varepsilon = 1\text{e-}07$.

The first set of risk levels we tested was $r_1 = (10\%, 50\%, 90\%)$. The estimated IRD for the logistic regression solution $\hat{\beta}_{LR}$ and the matching (non-degenerate) set of breakpoints $\hat{\tau}(\hat{\beta}_{LR}) = (-2.6918, 9.1698)$ was slightly above our set feasibility threshold $(\hat{\text{IRD}}_{r_1}(\hat{\beta}_{LR}) = 7.8776\text{e-}05)$. The second set of risk levels we tested was $r_2 = (20\%, 50\%, 80\%)$. For this set the solution $\hat{\beta}_{LR}$ provided by the logistic regression was clearly infeasible: on the one hand the interval associated to the 50% risk level was clearly degenerate $(\hat{\tau}_2(\hat{\beta}_{LR}) - \hat{\tau}_1(\hat{\beta}_{LR}) < 1.1\text{e-}06)$ and the estimated IRD was high above our set threshold $(\hat{\text{IRD}}_{r_2}(\hat{\beta}_{LR}) = 0.0014 > \varepsilon)$. We proceeded to construct a constrained maximum likelihood problem for both $r_1, r_2$ as described in section 5. For $r_1$ an unconstrained problem was sufficient, with IRD < 1e-07 and a non-degenerate $\hat{\tau}(\hat{\beta})$. For $r_2$, the unpenalized ordinal risk-group problem produced degenerate solutions, so we added the penalty function $Pen(\beta, \tau) = \left(\frac{\max|\tau_i - \tau_{i-1}|}{\beta^T(\hat{\mu}_1 - \hat{\mu}_0)} - 1\right)^2$, which was designed to balance between the distance between the breakpoints $\tau_1, \tau_2$ and the distance between estimated class means, and selected a penalty coefficient $\gamma = 10$. Since in many cases the optimization algorithm converged to a local minimum we randomly sampled 25,000 starting points for each set of risk categories we tested. The estimates for logistic regression (LR) and ordinal risk-group logistic regression (ORG-LR) for both $r_1, r_2$ are summarized in table 1.

|  | $r_1$ | | $r_2$ | |
|---|---|---|---|---|
|  | LR | ORG-LR | LR | ORG-LR |
| Intercept | -87.6641 | 12.3928 | -87.6641 | -3.1977 |
| Texture | 0.2864 | 0.1894 | 0.2864 | 0.1844 |
| log(Area) | 11.9706 | 0.8999 | 11.9706 | -0.3573 |
| Smoothness | 67.7342 | -25.2575 | 67.7342 | 4.0826 |
| log(Compactness) | -1.8107 | -3.9205 | -1.8107 | 0.3070 |
| log(Concave Points) | 3.6698 | 11.8320 | 3.6698 | 0.1162 |
| log(Symmetry) | 3.2556 | 0.7888 | 3.2556 | -0.3529 |
| log-Likelihood | -45.9909 | -94.8944 | -45.9909 | -311.019 |
| $\hat{\tau}_1$ | -2.6918 | -4.4029 | 3.117868 | -0.7088 |
| $\hat{\tau}_2$ | 9.1698 | 6.8524 | 3.117869 | 1.3376 |
| $P(Y = 1 \mid \Psi(\hat{\beta}, X) \leq \tau_1)$ | 9.2925% | 9.9751% | 16.7680% | 19.9857% |
| $P(Y = 1 \mid \Psi(\hat{\beta}, X) \in (\tau_1, \tau_2])$ | 50.3289% | 50.0130% | 51.5483% | 50.0116% |
| $P(Y = 1 \mid \Psi(\hat{\beta}, X) > \tau_2)$ | 89.5768% | 89.9870% | 78.7930% | 79.9948% |
| $\hat{IRD}$ | 7.88e-05 | 9.69e-08 | 0.0014 | 3.66e-08 |

Table 1: Maximum likelihood (ML) estimators of coefficients, log-likelihood, optimal breakpoints $\tau$, model predicted probabilities (assuming multivariate normal distribution) and IRD estimates for unconstrained logistic regression and ordinal risk-group logistic regression (ORG-LR) for $r_1 = (10\%, 50\%, 90\%)$, $r_2 = (20\%, 50\%, 80\%)$

The differences between the two methods can be further illustrated by looking at the distributions of the logit-transformed predicted probabilities $\{logit(\hat{Y}_i)\}_{i=1}^N$ of both methods. Figure 2 illustrates the logistic regression solution for $r_1$ (top graph) and the ordinal risk-group logistic regression solution for $r_1$ (bottom graph), and figure 3 illustrate the same results for $r_2$. A comparison of the two graphs in each figure shows that for both sets of risk levels, the ORG-LR solution compromises the quality of separation between the two classes in order to achieve feasibility, and in the more extreme case of $r_2$ reduces separation dramatically in order to avoid degenerate solutions.

In order to validate the results of our new method and compare them to the performance of the logistic regression solution we performed a cross validation study. We randomly divided the dataset into two groups: 90% of the patients were randomly sampled as a training set, from which a logistic regression model and an ordinal risk-group model were constructed, and the remaining 10% were used as a test set for the models. We repeated the process with 25,000 random samples and calculated the percentage of "Melignant" cases found in each predicted risk group for each of the models. The results of the implementation of the training models on the test sets for $r_1 = (10\%, 50\%, 90\%)$ were $(0.7173\%, 74.6978\%, 100\%)$ for logistic regression (IRD = 0.07962) and $(3.8804\%, 57.4549\%, 99.9329\%)$ for ordinal risk-group logistic regression (IRD = 0.01917). The cross validation results for $r_2 = (20\%, 50\%, 80\%)$ using ORG-LR were $(15.7349\%, 52.9507\%, 84.9715\%)$ (IRD = 0.0052). Since the logistic regression solution was degenerate for $r_2$ we were unable test it with cross validation.
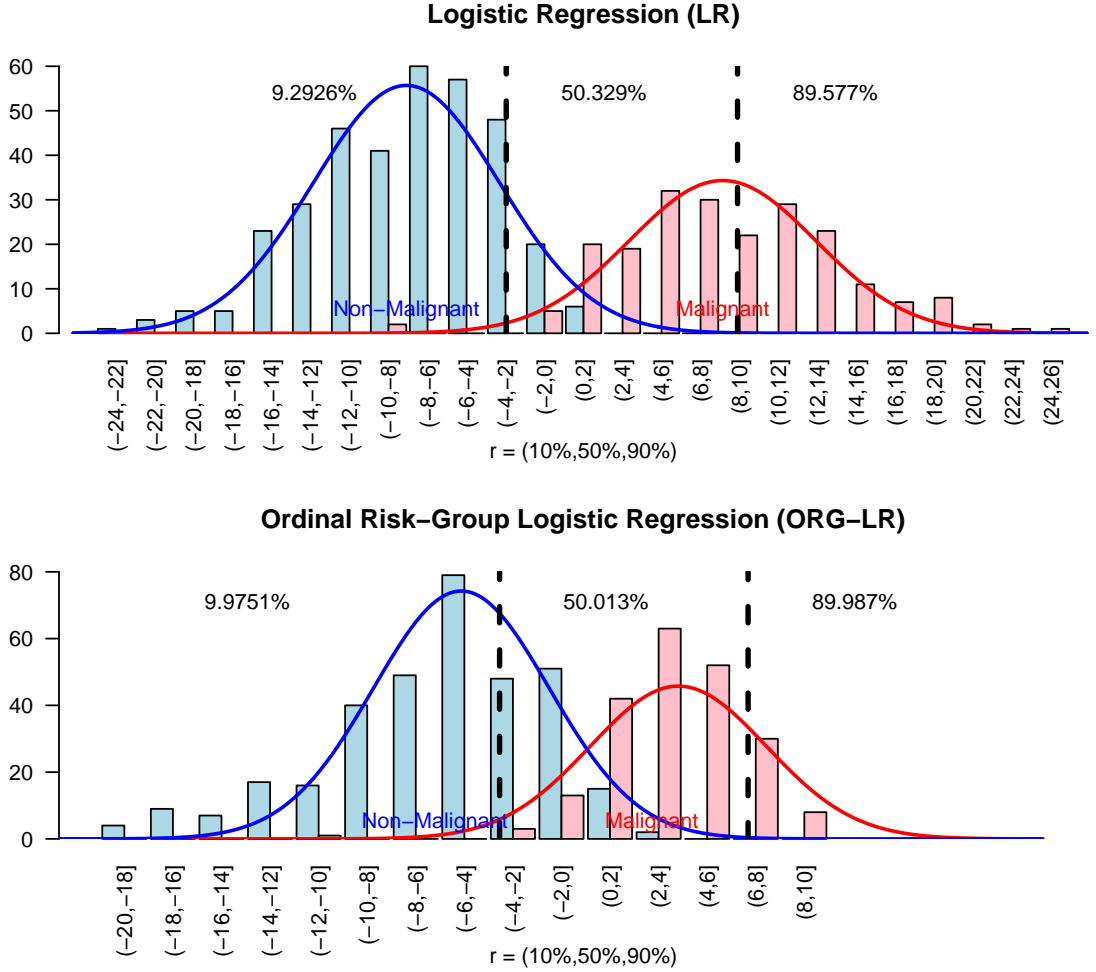
Figure 2: Logit-transformed predictions and separation between the malignant and non-malignant classes of logistic gegression (LR) (top graph) and the Ordinal Risk-Group Logistic Regression (ORG-LR) (bottom graph) for $r_1$. The black dotted lines mark the matching sets of breakpoints $\hat{\tau}(\hat{\beta}_{LR})$ and $\hat{\tau}(\hat{\beta}_{ORG-LR})$.

The comparison of the cross validation results from the two methods for $r_1$ shows that although both models did not perform wery well, the ORG-LR solution outperformed the logistic regression solution (IRD is approx. 4 times smaller), in spite of the fact that differences in IRD between the two models do not seem significant. We estimate that one of the reasons for the of the high absolute deviance in IRD of all models and risk levels we tested is that the data is not exactly normally distributed (as evident in figures 2,3). In addition, specifically for ORG-LR, we suspect that the generic algorithms we used for the ordinal risk-group maximum likelihood constrained optimization failed to converge to the real global minimum in some of the cross-validation iterations. Verifying this hypothesis would require either use of more specific optimization algorithms (for example by analytically calculating the derivatives of the IRD constraint) or by a very large scale simulation study that would require billions of iterations. Both approaches are outside the scope of this paper and we leave them for future studies.
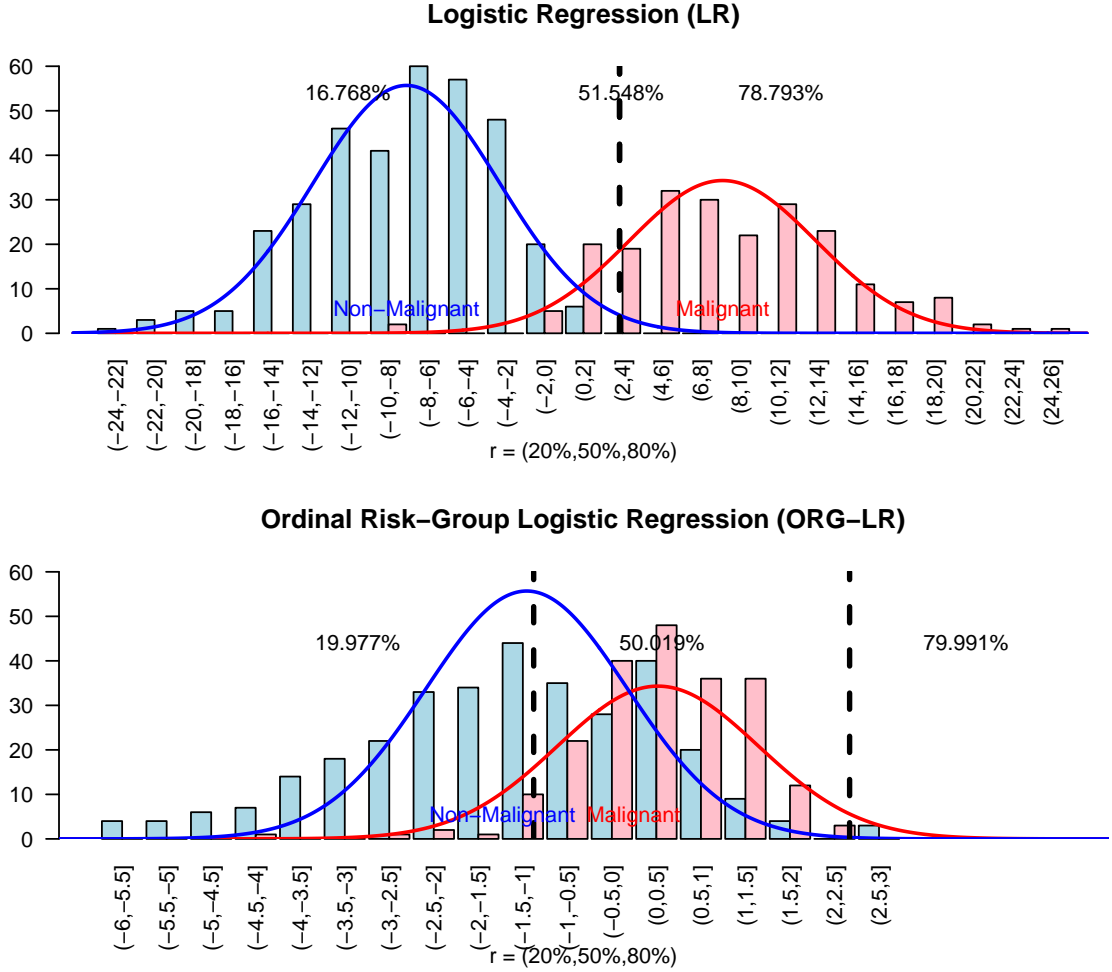
Figure 3: Logit-transformed predictions and separation between the malignant and non-malignant classes of logistic gegression (LR) (top graph) and the Ordinal Risk-Group Logistic Regression (ORG-LR) (bottom graph) for $r_2$. The black dotted lines mark the matching sets of breakpoints $\hat{\tau}(\hat{\beta}_{LR})$ and $\hat{\tau}(\hat{\beta}_{ORG-LR})$.

# 6  Conclusions

The exact estimation of the conditional risk function is an important part of practical and theoretical research. However the practical application of this information is very often in the form of a finite and small set of resulting actions. Although conditional risk quantiles provide valuable information, we ultimately want to know the risk associated with adjacent non-overlapping intervals in order to create distinct ordinal risk groups. As we have demonstrated in section 3.1, quantile regression is not useful for that purpose. Furthermore, section 3.2 demonstrates that the practice of dividing post-hoc the continuous estimate of conditional risk into intervals ignores the limitations introduced by the lower bounds on IRD and may produce sub-optimal or degenerate solutions.

Our formulation of the optimization problem, as presented in section 4, reflects our understanding that while the model's ability to separate the classes remains

the key issue, we must introduce both a new constraint and a penalty function in order to achieve two additional objectives: an accurate risk distribution and a usable partition scheme. While IRD represents an absolute measure of the model's quality and must be a constraint on the optimal solution, the "softer" requirement on minimal interval length should allow for flexibility in application. We believe that a penalty function enables better control and adaptation through the choice of function and the aversion parameter.

Finally, we wish to emphasize the implications of the most counter intuitive result of this paper - the existence of limitations on certain risk structures (the vector $r$) in the form of lower bounds on the error rate IRD (equation 2). Although most of the examples we described are linear or logistic models with Gaussian conditional distributions, the existence of lower bounds holds for any continuous risk estimator. A re-evaluation of the optimal properties of such estimators in the context of risk discretization is therefore required. We leave the specifics of applying these ideas to other classification methods as well as proofs of consistency and the construction of confidence intervals to future studies.

# References

[1] T. Amemiya. Qualitative Response Models: A Survey. *Journal of economic literature*, 19(4):1483–1536, 1981.

[2] T. Amemiya. *Advanced Econometrics*. Harvard Univ Pr, 1985.

[3] D.R Cox. *Research papers in statistics: Festschrift for J. Neyman*. Wiley, 1966.

[4] D.R Cox. *The analysis of binary data*. London: Chapman & Hall, 1969.

[5] J. S. Cramer. The origins of logistic regression. Tinbergen Institute Discussion Papers 02-119/4, Tinbergen Institute, 2002.

[6] R. M. Dudley. *Uniform central limit theorems*. Cambridge University Press, 1999.

[7] G. Elliott and R.P. Lieli. Predicting Binary Outcomes. *Manuscript, UCSD*, 2006.

[8] R.A. Fisher. The Use of Multiple Measurements in Taxonomic Problems. *Annals of Eugenics*, 7:179–188, 1936.

[9] Frank A. and Asuncion A. UCI machine learning repository, 2010.

[10] W.H. Greene and D.A. Hensher. *Modelling Ordered Choices: A Primer*. Cambridge Univ Pr, 2010.

[11] R. J. Hathaway. A constrained formulation of maximum-likelihood estimation for normal mixture distributions. *The Annals of Statistics*, 13(2):pp. 795–800, 1985.

[12] D. V. Hinkley. On the ratio of two correlated normal random variables. *Biometrika*, 56(3):635–639, 1969.

[13] R. Koenker and G.Jr. Bassett. Regression Quantiles. *Econometrica: journal of the Econometric Society*, 46(1):33–50, 1978.

[14] G.S. Maddala. *Limited-Dependent and Qualitative Variables in Econometrics*. Cambridge Univ Pr, 1986.

[15] C.F. Manski. Maximum Score Estimation of the Stochastic Utility Model of Choice. *Journal of Econometrics*, 3(3):205–228, 1975.

[16] C.F. Manski. Semiparametric Analysis of Discrete Response:: Asymptotic Properties of the Maximum Score Estimator. *Journal of Econometrics*, 27(3):313–333, 1985.

[17] C.F. Manski and T.S. Thompson. Operational Characteristics of Maximum Score Estimation* 1. *Journal of Econometrics*, 32(1):85–108, 1986.

[18] D. Martin. Early Warning of Bank Failure: A Logit Regression Approach. *Journal of Banking & Finance*, 1(3):249–276, 1977.

[19] D. McFadden. Conditional logit analysis of qualitative choice behavior. *Frontiers in Econometrics*, pages 105–142, 1973.

[20] J. A. Nelder and R. W. M. Wedderburn. Generalized linear models. *Journal of the Royal Statistical Society. Series A (General)*, 135(3):pp. 370–384, 1972.

[21] J.A. Ohlson. Financial Ratios and the Probabilistic Prediction of Bankruptcy. *Journal of Accounting Research*, 18(1):109–131, 1980.

[22] J. Powell. Estimation of Semiparametric Models. In R. Engle and D. McFadden, editors, *Handbook of Econometrics Vol. 4*. North Holland, 1994.

[23] R. L. Prentice and R. Pyke. Logistic disease incidence models and case-control studies. *Biometrika*, 66(3):403–411, 1979.

[24] R Development Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2011. ISBN 3-900051-07-0.

[25] W. N. Street, W. H. Wolberg, and O. L. Mangasarian. Nuclear feature extraction for breast tumor diagnosis. 1905(1):861–870, 1993.

[26] H. Theil. A multinomial extension of the linear logit model. *International Economic Review*, 10(3):251–59, 1969.

[27] K. Train. *Discrete Choice Methods with Simulation*. Cambridge Univ. Press, 2003.

[28] A. W. van der Vaart. *Asymptotic statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 1998.

# A    Appendix: Equivalence of strict monotonicity of CDF ratios over intervals and the Strict Monotone Likelihood Ratio Property

Let $X$ be a continuous $P$-dimensional random vector and $Y \in \{0,1\}$ a Bernoulli random variable representing class membership. Assume that the conditional distributions $X \mid Y = k$ $(k = 0,1)$ are also continuous and the conditional densities $f_{X|Y=k}$ are finite. Let $\Psi : \mathbb{R}^P \to \mathbb{R}$ be a finite continuous risk predictor. For a single observation $X$, the *likelihood ratio* of the risk predictor $\Psi$ between the two alternatives represented by $Y$ is defined as the ratio of the conditional densities:

$$\Lambda_\Psi(x) = \frac{P(\Psi(X) = x \mid Y = 1)}{P(\Psi(X) = x \mid Y = 0)} = \frac{f_{\Psi(X)|Y=1}(x)}{f_{\Psi(X)|Y=0}(x)}$$

The *Strictly Monotone Likelihood Ratio Property* (SMLRP) demands that for a given $X, Y, \Psi$ the ratio $\Lambda_\Psi(x)$ is a strictly monotone function of $x$. It is worth noting that by using Bayes theorem we can show that demanding SMLRP is equivalent to demanding strict monotonicity of $P(Y = 1 \mid \Psi(X) = x)$ in $x$:

$$P(Y = 1 \mid \Psi(X) = x) = \frac{p f_1(x)}{p f_1(x) + (1 - p) f_0(x)} = \frac{1}{1 + \frac{1-p}{p} \frac{1}{\Lambda_\Psi(x)}} \tag{44}$$

where $p = P(Y = 1)$ and $f_k(x) = f_{\Psi(X)|Y=k}(x)$ are the density functions of the conditional distributions. This equivalence means that in terms of conditional probability, SMLRP is equivalent to strict *pointwise* monotonicity in the condition, in contrast to the requirement of monotonicity over right-expanding intervals in section 3.2, which we defined as the strict monotonicity of $R(\beta, \tau)_i = P(Y = 1 \mid \Psi(X) \in (\tau_{i-1}, \tau_i])$ in $\tau_i$ for any $\tau_{i-1}$ while $\tau_i > \tau_{i-1}$.

**Theorem A.1.** *SMLRP $\Leftrightarrow \forall \tau_{i-1}, \forall \tau_i > \tau_{i-1}$ $R(\beta, \tau)_i$ is strictly increasing in $\tau_i$.*

*Proof.* Using our previous definition of $R(\beta, \tau)_i = P(Y = 1 \mid \Psi(X) \in (\tau_{i-1}, \tau_i])$ and Bayes theorem we can represent:

$$R(\Psi, \tau)_i = \frac{p(F_1(\tau_i) - F_1(\tau_{i-1}))}{p(F_1(\tau_i) - F_1(\tau_{i-1})) + (1-p)(F_0(\tau_i) - F_0(\tau_{i-1}))} = \frac{1}{1 + \frac{1-p}{p} \frac{1}{\gamma_\Psi(\tau_{i-1}, \tau_i)}} \tag{45}$$

where

$$\gamma_\Psi(\tau_{i-1}, \tau_i) = \frac{F_1(\tau_i) - F_1(\tau_{i-1})}{F_0(\tau_i) - F_0(\tau_{i-1})} \tag{46}$$

and $F_k(x) = F_{\Psi(X)|Y=k}(x)$ are the cumulative distribution functions (CDF) of the conditional distributions. The strict monotonicity of $R_i(\Psi, \tau)$ in $\tau_i$ for any $\tau_{i-1} < \tau_i$ is therefore equivalent to the strict monotonicity of $\gamma(c, x)$ in $x$ for any $c$, $x > c$.

In addition, for two positive, finite, strictly increasing and once differentiable functions $g, h$ the following equivalence holds:

$$\frac{g(x)}{h(x)} \text{ is strictly increasing} \Leftrightarrow \frac{g'(x)h(x) - g(x)h'(x)}{h^2(x)} > 0 \Leftrightarrow \frac{g(x)}{h(x)} < \frac{g'(x)}{h'(x)} \tag{47}$$

Since $F_0, F_1$ meet these requirements, then by (45) the strict monotonicity of both $\gamma_\Psi(c,x)$ and $R(\Psi,\tau)_i$ is equivalent to following condition:

$$\gamma_\Psi(c,x) = \frac{F_1(x) - F_1(c)}{F_0(x) - F_0(c)} < \frac{\frac{d}{dx}(F_1(x) - F_1(c))}{\frac{d}{dx}(F_0(x) - F_0(c))} = \frac{f_1(x)}{f_0(x)} = \Lambda_\Psi(x) \quad \forall c < x \quad (48)$$

It is therefore sufficient to show that under the above assumptions of continuity and finiteness that the following equivalence holds:

$$SMLRP \iff \gamma_\Psi(c,x) < \Lambda_\Psi(x) \quad \forall c, c < x \quad (49)$$

*Step 1: $SMLRP \implies \gamma_\Psi(c,x) < \Lambda_\Psi(x) \quad \forall c < x$*
Under SMLRP:

$$\forall x_1 > x_0 \quad \frac{f_1(x_1)}{f_0(x_1)} > \frac{f_1(x_0)}{f_0(x_0)} \iff f_1(x_1)f_0(x_0) > f_1(x_0)f_0(x_1) \quad (50)$$

The equivalence holds since $f_0, f_1$ are strictly positive, continuous and finite. Integrating on $x_0$ over the interval $[c, x_1]$ we have:

$$\int_c^{x_1} f_1(x_1)f_0(x_0)dx_0 > \int_c^{x_1} f_1(x_0)f_0(x_1)dx_0$$
$$\iff f_1(x_1)(F_0(x_1) - F_0(c)) > f_0(x_1)(F_1(x_1) - F_1(c)) \iff \gamma_\Psi(c,x_1) < \Lambda_\Psi(x_1) \quad (51)$$

and this holds for any $x_1 \in \mathbb{R}$. In addition setting $c = -\infty$ when integrating maintains the strict inequalities of (51), and therefore SMLRP also ensures strict monotonicity of $\gamma_\Psi(-\infty, x) = F_1(x)/F_0(x)$ and the equivalent monotonicity of $P(Y = 1 \mid \Psi(X) < x)$ in $x$.

*Step 2: $\forall c, c < x \quad \gamma_\Psi(c,x) < \Lambda_\Psi(x) \implies SMLRP$*

Under our assumptions:

$$\forall x_1 > x_0 > c \quad \gamma(c,x_1) = \frac{\int_c^{x_1} f_1(x)dx}{\int_c^{x_1} f_0(x)dx} > \frac{\int_c^{x_0} f_1(x)dx}{\int_c^{x_0} f_0(x)dx} = \gamma(c,x_0) \quad (52)$$

Assuming all functions are continuous and finite we can take $c \to x_0$:

$$\gamma(x_0, x_1) = \frac{\int_{x_0}^{x_1} f_1(x)dx}{\int_{x_0}^{x_1} f_0(x)dx} > \frac{f_1(x_0)}{f_0(x_0)} \quad (53)$$

This inequality is strict since for any $x_1 > x_0$ there exists $\varepsilon = \frac{x_1 - x_0}{2} > 0$ such that:

$$\gamma(x_0, x_1) = \frac{\int_{x_0}^{x_1} f_1(x)dx}{\int_{x_0}^{x_1} f_0(x)dx} > \frac{\int_{x_0}^{x_0+\varepsilon} f_1(x)dx}{\int_{x_0}^{x_0+\varepsilon} f_0(x)dx} = \gamma(x_0, x_0 + \varepsilon) \geq \frac{f_1(x_0)}{f_0(x_0)} \quad (54)$$

On the other hand taking $c \to x_1$ (using the same considerations and utilizing the fact that for $b > a$, $\int_b^a f(x)dx = -\int_a^b f(x)dx$) and combining with (53) we have:

$$\frac{f_1(x_1)}{f_0(x_1)} > \frac{\int_{x_0}^{x_1} f_1(x)dx}{\int_{x_0}^{x_1} f_0(x)dx} = \gamma(x_0, x_1) > \frac{f_1(x_0)}{f_0(x_0)} \quad (55)$$

$\square$

# B  Appendix: Uniqueness of $\tau$ under the Strict Monotone Likelihood Ratio Property

The risk estimation methods mentioned in this paper typically deal only with the optimal estimation of $\Psi$ (and the breakpoints $\tau$ are defined post-hoc). Introducing the set of breakpoints $\tau$ as an integral part of the definition of IRD increases in the number of parameters that must be estimated simultaneously, resulting in a more complicated parameter space (for example we require $\tau_{i-1} < \tau_i$). Although the increase in the number of estimated parameters should not be significant (in practical scenarios we expect $T \leq 10$) the result nonetheless would be longer running times for the optimization algorithms. Before we proceed any further it would be useful to identify sufficient conditions for the uniqueness of $\tau$ for a given $\Psi$:

**Lemma B.1.** *If for a given $\Psi$ the likelihood ratio $\Lambda_\Psi(x) = \frac{f_{\Psi(X)|Y=1}(x)}{f_{\Psi(X)|Y=0}(x)}$ satisfies the strict monotone likelihood ratio property (SMLRP), then if there exists $\tau_\Psi$ such that $IRD_r(\Psi, \tau_\Psi) = 0$ it is unique.*

*Proof.* If $\Lambda_\Psi(x)$ satisfies SMLRP, then by theorem A.1 (appendix A) $R$ is strictly monotone is $\tau_i$. The rest is by induction: strict monotonicity of $R_1$ means that if there exists $\tau_1$ which satisfies $R(\Psi, \tau)_1 = r_1$, then it is unique. Fixing $\tau_{i-1}$, if (11) holds (meaning that $\tau_i$ is "feasible"), then again by strict monotonicity, if there exists $\tau_i$ that satisfies $R(\Psi, \tau)_i = r_i$, then it is unique.

Therefore if (11) holds for all $i$ then only a single $\tau$ satisfies $IRD_r(\Psi, \tau) = 0$. $\square$

**Corollary B.2.** *Under SMLRP we can denote $\tau = \tau(\Psi)$ and define IRD using $\Psi$ alone:*

$$R_i(\Psi) = P(Y = 1 \mid \Psi(X) \in (\tau_{i-1}(\Psi), \tau_i(\Psi)]), \quad IRD_r(\Psi) = \|R(\Psi) - r\| \quad (56)$$