A Personalized Boosting Screening Method and Its Application to Sepsis

^aSchool of Industrial & Systems Engineering, Georgia Tech, Atlanta, GA; ^bRegenstrief Center for Healthcare Engineering, Purdue University, West Lafayette, IN; ^cCritical Care Medicine, Northeast Georgia Medical Center, Gainesville, GA

Motivation and Objective

- The natural mortality rate for **sepsis** is between 25% and 50%. Nearly half of patients who die in hospitals are septic
- In 2016, a task force committee recommended screening for sepsis by quick Sepsis-related Organ Failure Assessment (qSOFA), which uses the constant thresholds in decision making.



• Our Objective: Develop a personalized sepsis screening method that depends on a patient's baseline characteristics such as age, sex, admission location, etc.

Problem Formulation

- Data: (Y_i, X_i, u_i) , for $i = 1, \dots, N$, where $Y_i \in \{-1, 1\}$ is the binary outcome, X_i is the biomarker (e.g., blood pressure, respiratory rate, etc.), and $u_i \in R^q$ are baseline characteristics. - Classification Rule: Predict Y_i by

if $X_i \geq c_i(\boldsymbol{u}_i)$ $\hat{Y}_i =$ otherwise $= \operatorname{sign}(X_i - c_i(\boldsymbol{u}_i)).$

Here we assume

$$c_i(\boldsymbol{u}_i) = \boldsymbol{u}_i^T oldsymbol{eta}$$

- for some unknown parameters $\beta = (\beta_0, \beta_1, \cdots, \beta_q)^T$. • Question: Estimate β so as to minimize misclassification rate. • The function sign(f) is not continuous.
- The 0-1 loss function $I(Y \neq sign(f))$ is non-smoothing.
- The consequences of misclassifying sepsis and non-sepsis patients are different.

Existing Methods

qSOFA considers the constant threshold, i.e., $c_i(\boldsymbol{u}_i) \equiv c$. Several approaches to find suitable constant thresholds: • minP Approach: Maximizing the standard chi-square statistic • Youden Index: Maximizing the sum of sensitivity and specificity Closest-to-(0,1) Criterion: The "optimal" threshold is defined as the point on the ROC curve closest to (0,1)

Chen Feng^a, Paul Griffin^b, Shravan Kethireddy^c, and Yajun Mei^a

Our Proposed Method: Personalized Threshold

Key ideas in our proposed method

- We define the threshold $c_i(\boldsymbol{u}_i)$ as a function of the individual subjects' characteristics \boldsymbol{u}_i
- We borrow the idea of **boosting** to replace the 0-1 loss function $I(Y \neq sign(f))$ by the exponential loss exp(-Yf)
- Introducing two different weights, w_+ and w_- , depending on whether $Y_i = +1$ or -1, in order to take into account the different consequences of misclassification. **Parameter Estimation as Optimization Problem**

$J(\boldsymbol{\beta}) = \frac{1}{N} \sum_{i=1}^{N} $	$w_+ \cdot e^{-Y_i f_i} \cdot I(Y_i)$
$= \frac{1}{N} \sum_{i=1}^{N} ($	$\left(w_{Y_i} \cdot e^{-Y_i f_i}\right) = \frac{1}{N}$

where $f_i \triangleq f(X_i, \boldsymbol{u}_i) = X_i - c_i(\boldsymbol{u}_i)$ and $w_{Y_i} \triangleq (w_+ \cdot (Y_i + 1) + w_-)$ $(1 - Y_i))/2.$

Computational algorithm: Gradient descent

Algorithm 1 Gradient Descent Using Weighted Exponential Loss Require: $Y, X, U, N, w_{-} > 0, w_{+} > 0, \alpha, T$ Le Initialization: $\beta_i \leftarrow 0 \ \forall i \in \{0, 1, 2, ..., q\}, \ W = (w_+(Y+1) + w_-(1-Y))/2$ 2: for all t = 1, 2, ..., T do 3: $\boldsymbol{f} = \boldsymbol{X} - \boldsymbol{U}^T \boldsymbol{eta}$ $L = \exp(-Y * f)$ {*: Element-wise product} Forward Propagation 5: $\boldsymbol{J} = \frac{1}{N} \boldsymbol{L}^T \boldsymbol{W}$ 6: $df = -\frac{1}{N}W * Y * \exp(-Y * f)$ **Backward** Propagation 7: $d\boldsymbol{\beta} = -\boldsymbol{U}d\boldsymbol{f}$ 8: $\boldsymbol{\beta} \leftarrow \boldsymbol{\beta} - \alpha d\boldsymbol{\beta}$ 9: end for 10: $\hat{c} = \boldsymbol{U}^T \boldsymbol{\beta}$

Proposition: The objective function $J(\beta)$ is convex with respect to β , and thus the gradient descent algorithm converges to the global optimum if the learning rate α is small enough and the optimization steps T is long enough.

Medical Information Mart for Intensive Care III (MIMIC-III) database (version 1.4) **Study Population:** 3,771 sepsis patients (Y = 1); 4,000 non-sepsis patients (Y = -1)**qSOFA variables (**X**)**: respiratory rate (RR), systolic blood pressure (sysBP), and GCS scores (we keep the constant cutoff 15 for GCS score, since GCS score is a discrete variable) **Baseline characteristics** (u's): age, gender, admission location, admission type, ethnicity, insurance, and marital status

We propose to estimate the (q+1)-dimensional parameter β by minimizing the training error under the weighted exponential loss function:

$$= 1) + w_{-} \cdot e^{-Y_{i}f_{i}} \cdot I(Y_{i} = -1)),$$

$$\sum_{i=1}^{N} \left(w_{Y_{i}} \cdot e^{-Y_{i}\left(X_{i} - \boldsymbol{\beta}^{T}\boldsymbol{u}_{i}\right)} \right)$$

$$(1)$$

$$(1)$$

Data Set

Comparison to qSOFA thresholds and existing methods

The parameters T = 30000, $\alpha = 0.001$, $w_+ = 1$, and $w_- = 1$ were selected based on a grid search to maximize the prediction accuracy. **Overall Accuracy**



Perso Log

Intensive Care Unit

Application to Sepsis Screening

C) Compared to Machine Learning Techniques

Methods	Overall Accuracy	Sensitivity	Specificity
nalized qSOFA	0.7093	0.6668	0.7493
istic Regression	0.7489	0.7428	0.7548
AdaBoosting	0.7456	0.7359	0.7549
		-	

• Machine Learning methods are black boxes and difficult for clinicians and nurses to implement for real time monitoring in

Contact information