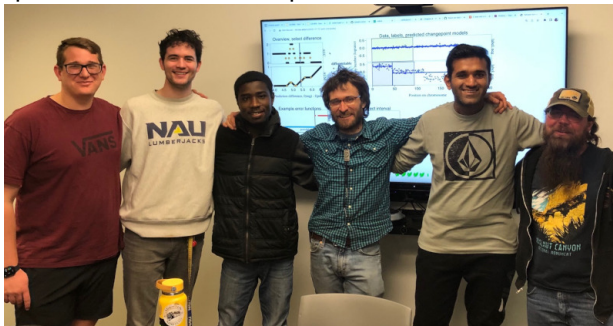


Sort-based gradient descent learning algorithms for ROC curve optimization

Toby Dylan Hocking — toby.dylan.hocking@usherbrooke.ca
joint work with former master student Jadon Fowler
Learning Algorithms, Statistical Software, Optimization
(LASSO Lab) — <https://lassolab.org>
Département d'Informatique, Université de Sherbrooke



Problem Setting 1: ROC curves for evaluating supervised binary classification algorithms

Proposed AUM loss, a differentiable surrogate for ROC curve optimization (JMLR'23)

Problem setting 2: ROC curves for evaluating supervised changepoint algorithms

Proposed complete line search algorithm for optimizing AUM and AUC (in progress)

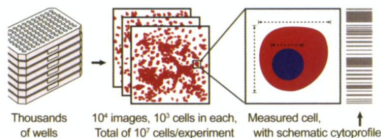
Discussion and Conclusions

Problem: supervised binary classification

- ▶ Given pairs of inputs $\mathbf{x} \in \mathbb{R}^p$ and outputs $y \in \{0, 1\}$ can we learn a score $f(\mathbf{x}) \in \mathbb{R}$, predict $y = 1$ when $f(\mathbf{x}) > 0$?
- ▶ Example: email, \mathbf{x} = bag of words, y = spam or not.
- ▶ Example: images. Jones *et al.* PNAS 2009.

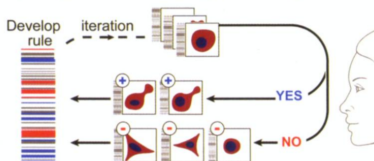
A Automated Cell Image Processing

Cytoprofile of 500+ features measured for each cell



B Iterative Machine Learning

System presents cells to biologist for scoring, in batches

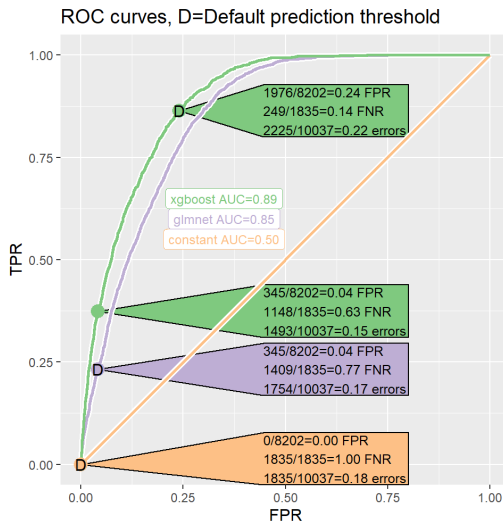


Most algorithms (SVM, Logistic regression, etc) minimize a differentiable surrogate of zero-one loss = sum of:

False positives: $f(\mathbf{x}) > 0$ but $y = 0$ (predict budding, but cell is not).

False negatives: $f(\mathbf{x}) < 0$ but $y = 1$ (predict not budding, but cell is).

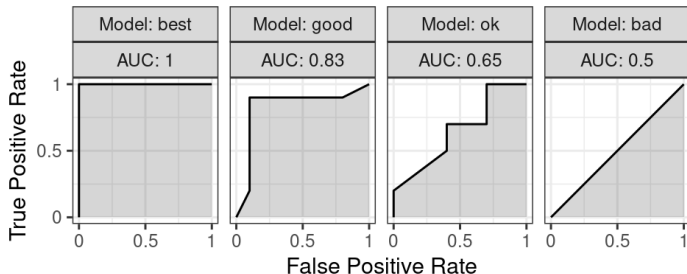
ROC curves for evaluation, especially useful with imbalance



- ▶ At default FPR=24% (D), glmnet has fewer errors.
- ▶ At FPR=4%, xgboost has fewer errors.

Receiver Operating Characteristic (ROC) Curves

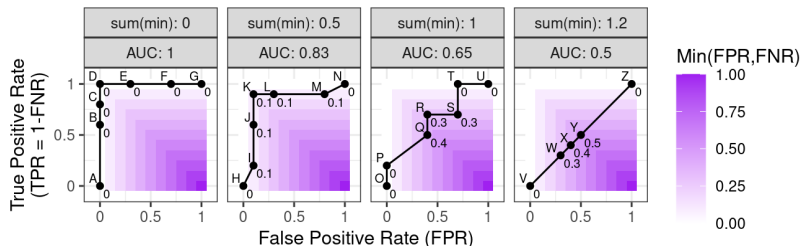
- ▶ Classic evaluation method from the signal processing literature (Egan and Egan, 1975).
- ▶ ROC curve of learned f is plot of True Positive Rate vs False Positive Rate: each point on the ROC curve is a different constant $c \in \mathbb{R}$ added to the predicted values: $f(\mathbf{x}) + c$.
- ▶ $c = \infty$ means always predict positive label (FPR=TPR=1).
- ▶ $c = -\infty$ means always predict negative label (FPR=TPR=0).
- ▶ Best classifier has a point near upper left (TPR=1, FPR=0), with large Area Under the Curve (AUC).



Research question and new idea

Can we learn a binary classification function f which directly optimizes the ROC curve?

- ▶ Most algorithms involve minimizing a differentiable surrogate of the zero-one loss, which is not the same.
- ▶ The Area Under the ROC Curve (AUC) is piecewise constant (gradient zero almost everywhere), so can not be used with gradient descent algorithms.
- ▶ We proposed (Hocking, Hillman 2023) to encourage points to be in the upper left of ROC space, using a loss function which is a differentiable surrogate of the sum of $\min(\text{FPR}, \text{FNR})$.



Problem Setting 1: ROC curves for evaluating supervised binary classification algorithms

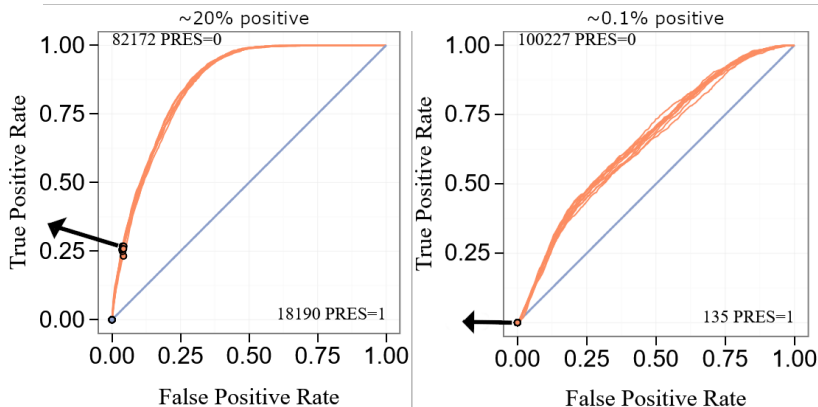
Proposed AUM loss, a differentiable surrogate for ROC curve optimization (JMLR'23)

Problem setting 2: ROC curves for evaluating supervised changepoint algorithms

Proposed complete line search algorithm for optimizing AUM and AUC (in progress)

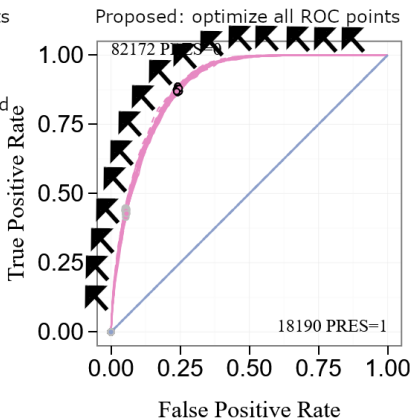
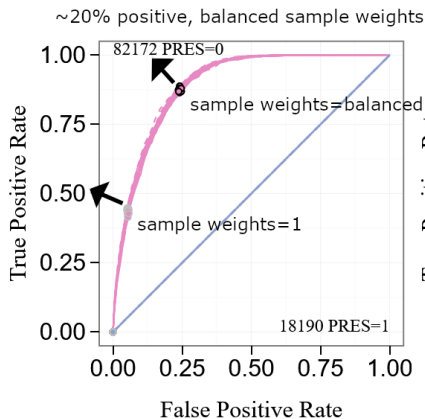
Discussion and Conclusions

Gradients of sample-based loss are influenced by imbalance



- ▶ Left: some imbalance, 20% positive labels, gradient 4x stronger along X axis / False Positive Rate.
- ▶ Right: large imbalance, 0.1% positive labels, gradient 1000x stronger along X axis / False Positive Rate. (True Positive / Y axis gradients essentially ignored)

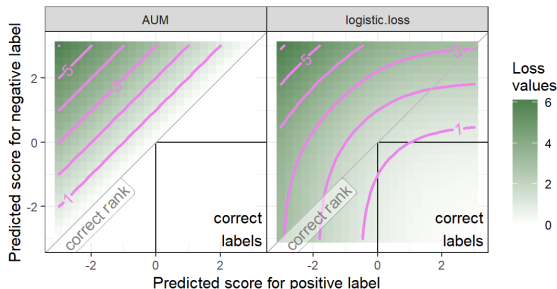
Gradients using balanced sample weights, proposed loss



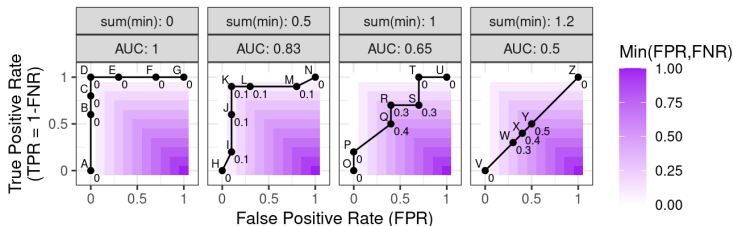
- ▶ Left: gradient 4x stronger along X axis for sample weights=1. Balanced sample weights mean equal influence for gradients along both axes, based on the current prediction threshold.
- ▶ Right: proposed method computes gradients based on all ROC points, not just the current prediction threshold.

Comparing proposed loss with baselines

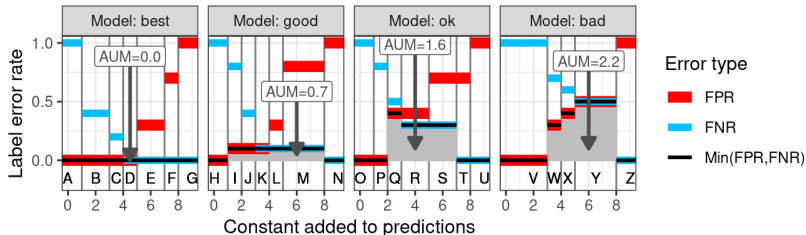
- ▶ Classic baselines: hinge and logistic loss, sum over samples, $\ell[yf(x)]$.
- ▶ Bamber (1975) proved ROC-AUC relation to Mann-Whitney U statistic (double sum over all pairs of positive and negative samples).
- ▶ Recently: SVM^{struct} (Joachims 2005), X-risk (Yang 2022), All Pairs Squared Hinge (Rust and Hocking 2023), sum loss over pairs of positive and negative samples, $\ell[f(x^+) - f(x^-)]$.
- ▶ Proposed: sort-based AUM loss (sum over points on ROC curve).
- ▶ Figure below: loss for two samples: one positive, one negative.



Large AUC \approx small Area Under Min(FP,FN) (AUM)

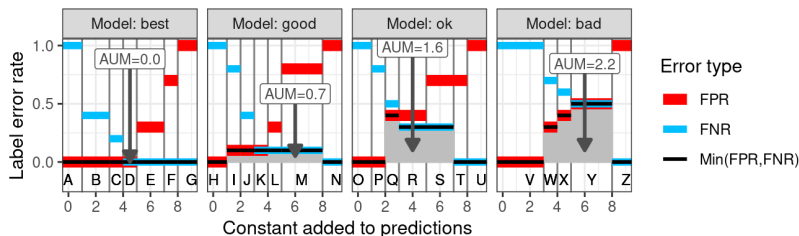


Above: purple heat map = numbers near dots = distance to top or left
= same as black min error rate functions below.



Hocking, Hillman, *Journal of Machine Learning Research* (2023).

Computing Sum of Min (SM)

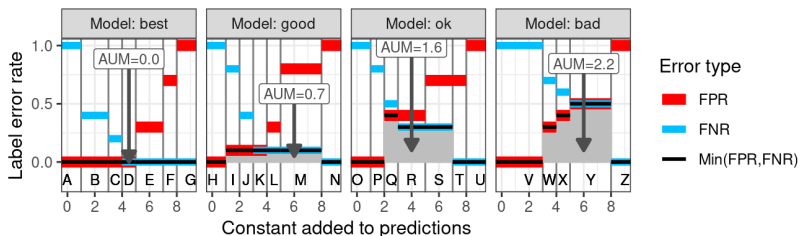


- ▶ For N samples, there are $\leq N + 1$ points on the ROC curve,
- ▶ with sorted thresholds $T_1 \leq \dots \leq T_N \in \mathbb{R}$,
- ▶ and corresponding min error values M_1, \dots, M_N .
- ▶ Then if I is the indicator function, we can write the sum of the min (SM), over all ROC points, as:

$$SM = \sum_{i=2}^N I[T_i \neq T_{i-1}] M_i = \sum_{i: T_i \neq T_{i-1}} M_i.$$

(\neq required: a tie $T_i = T_{i-1}$ deletes a point from the ROC curve)

Computing Area Under Min (AUM)



The AUM can be interpreted as an L1 relaxation of SM,

$$SM = \sum_{i=2}^N I[T_i \neq T_{i-1}] M_i = \sum_{i: T_i \neq T_{i-1}} M_i.$$

$$AUM = \sum_{i=2}^N [T_i - T_{i-1}] M_i.$$

AUM is therefore a surrogate loss for ROC optimization, and it is differentiable almost everywhere \Rightarrow gradient descent learning!

ROC curve pytorch code

```
def ROC_curve(pred_tensor, label_tensor):  
    sorted_indices = torch.argsort(-pred_tensor)  
    ...  
    return { # a dictionary of torch tensors  
        "FPR":FPR,  
        "FNR":FNR,  
        "TPR":1 - FNR,  
        "min(FPR,FNR)":torch.minimum(FPR, FNR),  
        "min_constant":torch.cat([  
            torch.tensor([-torch.inf]), uniq_thresh]),  
        "max_constant":torch.cat([  
            uniq_thresh, torch.tensor([torch.inf])])  
    }
```

<https://tdhock.github.io/blog/2024/torch-roc-aum/>

AUC and AUM pytorch code uses argsort

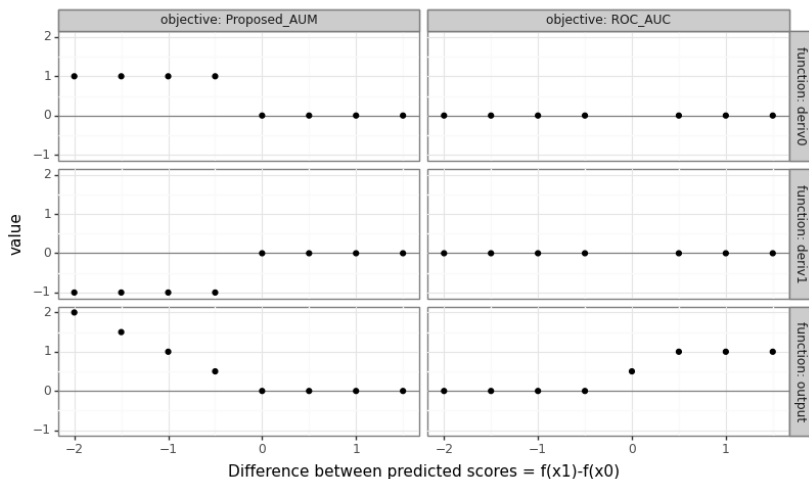
- ▶ ROC AUC and proposed AUM are both implemented by first computing the ROC curve.

```
def ROC_AUC(pred_tensor, label_tensor):  
    roc = ROC_curve(pred_tensor, label_tensor)  
    FPR_diff = roc["FPR"][1:] - roc["FPR"][:-1]  
    TPR_sum = roc["TPR"][1:] + roc["TPR"][:-1]  
    return torch.sum(FPR_diff * TPR_sum / 2.0)  
  
def Proposed_AUM(pred_tensor, label_tensor):  
    roc = ROC_curve(pred_tensor, label_tensor)  
    min_FPR_FNR = roc["min(FPR,FNR)"][1:-1]  
    constant_diff = roc["min_constant"][1:].diff()  
    return torch.sum(min_FPR_FNR * constant_diff)
```

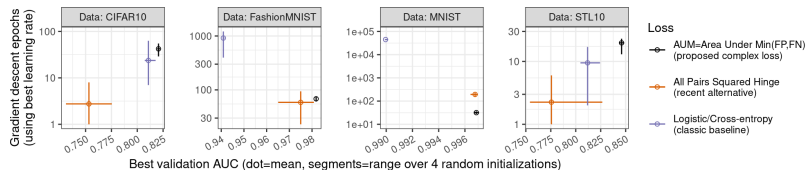
<https://tdhock.github.io/blog/2024/torch-roc-aum/>

AUM pytorch code, auto-grad demo

- ▶ Assume two samples, $(x_0, y_0 = 0), (x_1, y_1 = 1)$,
- ▶ Plot objective and gradient with respect to predicted scores.



AUM gradient descent increases validation AUC, four image classification data sets



- ▶ Unbalanced binary classification: 10% negative, 90% positive.
- ▶ Gradient descent with constant step size, best of 10^{-4} to 10^5 .
- ▶ Full gradient (batch size = number of samples).
- ▶ Linear model, max iterations = 100,000.
- ▶ Max Validation AUC comparable or better than baselines.
- ▶ Number of epochs comparable to baselines.
- ▶ Time per epoch is $O(N \log N)$ (sort), small log factor larger than standard logistic/cross-entropy loss, $O(N)$.

Discussion

- ▶ Classic baselines are simple (sum over samples).
- ▶ Proposed AUM loss similar to recent losses that sum over all pairs of positive and negative examples (both can be implemented by sorting predicted scores); proposed AUM loss uses a different/novel relaxation.
- ▶ Proposed AUM loss can be used as a drop-in replacement for logistic/binary cross-entropy loss,
<https://tdhock.github.io/blog/2024/torch-roc-aum/>
- ▶ Best use with stochastic gradient algorithms? At least one positive and one negative example is required in each batch.
- ▶ Margin-based algorithms like SVM?
- ▶ Works well in binary classification, how to adapt to multi-class setting, or other problems such as ranking/information retrieval? (next: change-point detection)

Problem Setting 1: ROC curves for evaluating supervised binary classification algorithms

Proposed AUM loss, a differentiable surrogate for ROC curve optimization (JMLR'23)

Problem setting 2: ROC curves for evaluating supervised changepoint algorithms

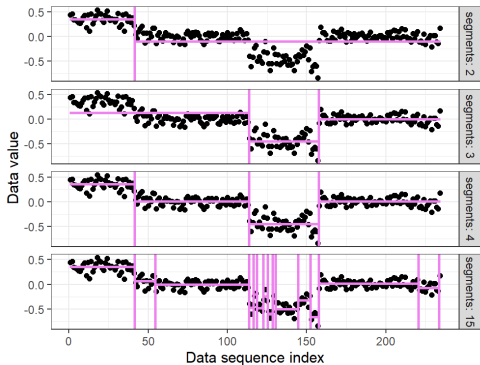
Proposed complete line search algorithm for optimizing AUM and AUC (in progress)

Discussion and Conclusions

Problem: unsupervised changepoint detection

- ▶ Data sequence z_1, \dots, z_T at T points over time/space.
- ▶ Ex: DNA copy number data for cancer diagnosis, $z_t \in \mathbb{R}$.
- ▶ The penalized changepoint problem (Maidstone *et al.* 2017)

$$\arg \min_{u_1, \dots, u_T \in \mathbb{R}} \sum_{t=1}^T (u_t - z_t)^2 + \lambda \sum_{t=2}^T I[u_{t-1} \neq u_t].$$

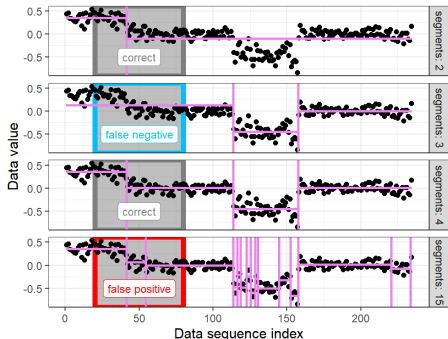


Larger penalty λ
results in fewer
changes/segments.

Smaller penalty
 λ results in more
changes/segments.

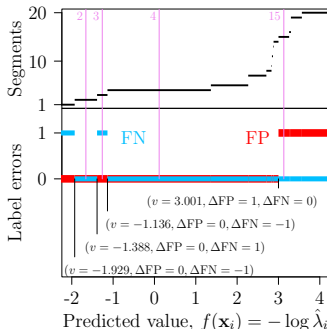
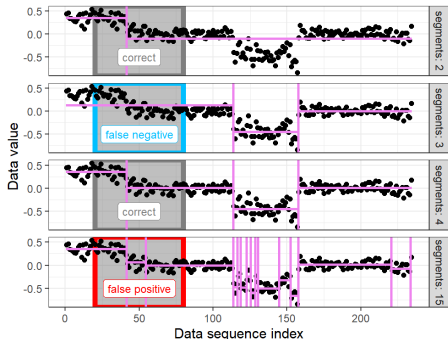
Problem: weakly supervised changepoint detection

- ▶ First described by Hocking *et al.* ICML 2013.
- ▶ We are given a data sequence \mathbf{z} with labeled regions L .



Problem: weakly supervised changepoint detection

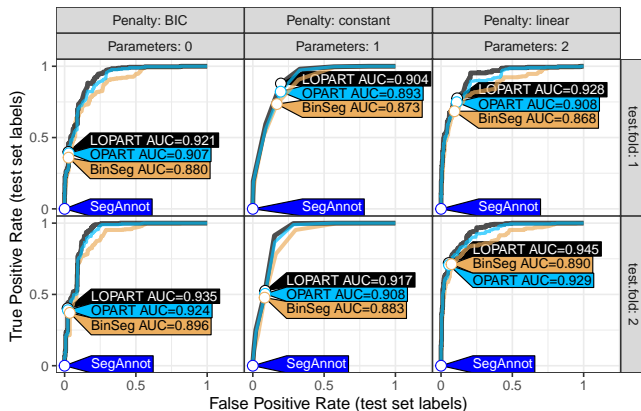
- ▶ First described by Hocking *et al.* ICML 2013.
- ▶ We are given a data sequence \mathbf{z} with labeled regions L .



We compute features $\mathbf{x} = \phi(\mathbf{z}) \in \mathbf{R}^p$ and want to learn a function $f(\mathbf{x}) = -\log \lambda \in \mathbf{R}$ that minimizes label error (sum of false positives and false negatives), or maximizes AUC, Hocking, Hillman, *Journal of Machine Learning Research* (2023).

Comparing changepoint algorithms using ROC curves

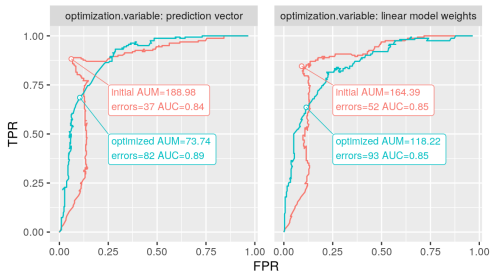
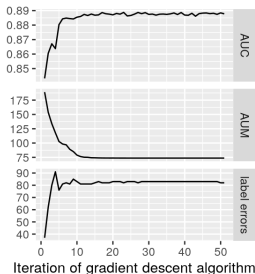
Hocking TD, Srivastava A. Labeled Optimal Partitioning.
Computational Statistics (2022).



LOPART algorithm (R package LOPART) has consistently larger test AUC than previous algorithms.

AUM gradient descent results in increased train AUC for a real changepoint problem

Hillman, Hocking, *Journal of Machine Learning Research* (2023).



- ▶ Left/middle: changepoint problem initialized to prediction vector with min label errors, gradient descent on prediction vector.
- ▶ Right: linear model initialized by minimizing regularized convex loss (surrogate for label error, Hocking *et al.* ICML 2013), gradient descent on weight vector.

Problem Setting 1: ROC curves for evaluating supervised binary classification algorithms

Proposed AUM loss, a differentiable surrogate for ROC curve optimization (JMLR'23)

Problem setting 2: ROC curves for evaluating supervised changepoint algorithms

Proposed complete line search algorithm for optimizing AUM and AUC (in progress)

Discussion and Conclusions

Using thresholds to compute AUM

Hillman and Hocking, JMLR 2023, showed that in a data set with N samples and $B \geq N$ breakpoints in label error functions, the AUM for the **current predictions** can be computed as a function of $T_1 < \dots < T_B$, thresholds when min label error M_b changes,

$$\text{AUM} = \sum_{b=2}^B [T_b - T_{b-1}] M_b$$

Proposed: when learning a linear model, $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$, we can update the weights \mathbf{w} using AUM gradient descent. We compute an exact representation of thresholds $T_b(s)$ and min error $M_b(s)$ **as a function of step size s** , which results in a complete piecewise linear representation of

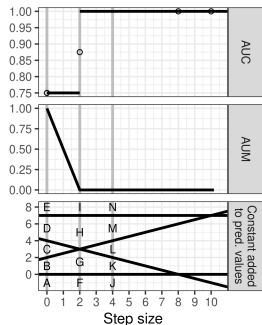
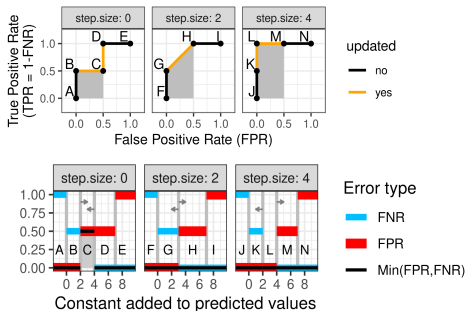
$$\text{AUM}(s) = \sum_{b=2}^B [T_b(s) - T_{b-1}(s)] M_b(s).$$

Simple example, proposed line search with 4 binary labels

Theorem: when learning a linear model, $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$,

- ▶ AUC is piecewise constant, and
- ▶ AUM is piecewise linear,

as a function of step size in gradient descent.



ROC curves and error functions

Proposed line search

Letters show correspondence between points on the ROC curve, and intervals of constants added to predicted values.

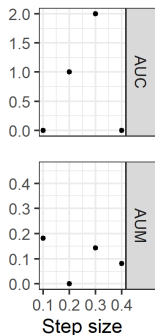
Change-point example, comparison with grid search

Theorem: when learning a linear model, $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$,

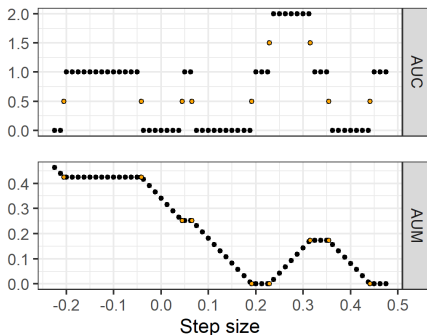
- ▶ AUC is piecewise constant, and
- ▶ AUM is piecewise linear,

as a function of step size in gradient descent.

Four steps



Many step sizes considered in grid search



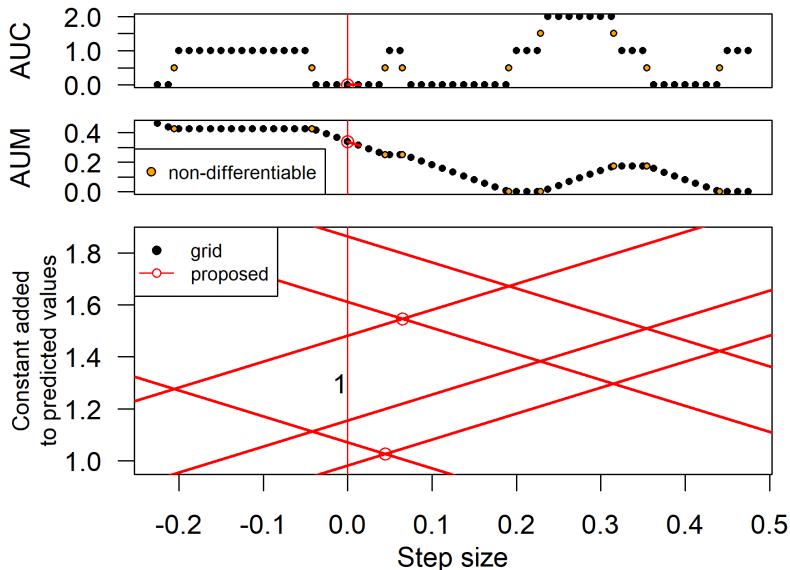
differentiable

- FALSE
- TRUE

Proposed line search algorithm computes updates when there are possible changes in slope of AUM / values of AUC (orange dots).

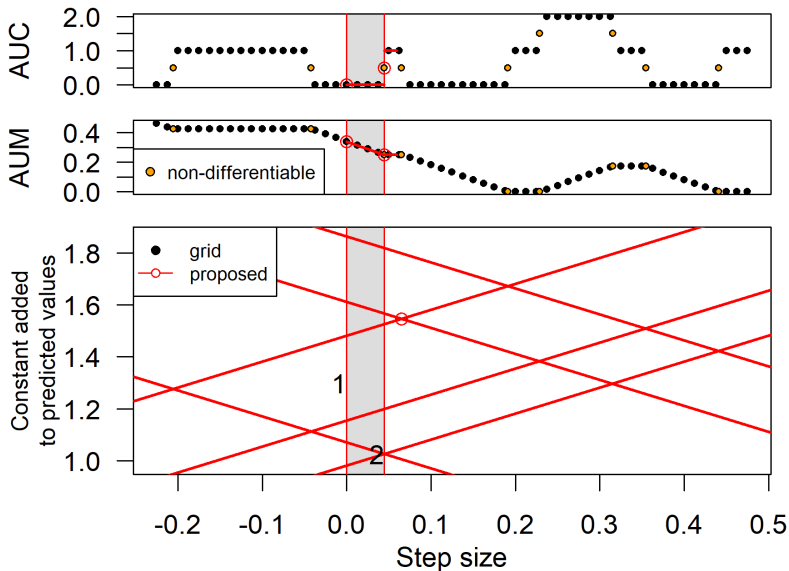
Proposed complete AUC/AUM line search, iteration 1

AUC/AUM values known only at red vertical line.



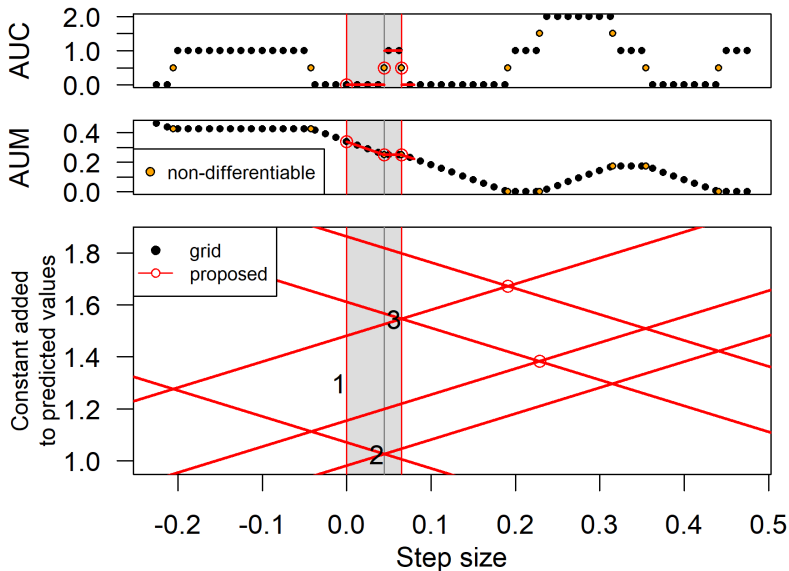
Proposed complete AUC/AUM line search, iteration 2

AUC/AUM values completely known within shaded grey region.



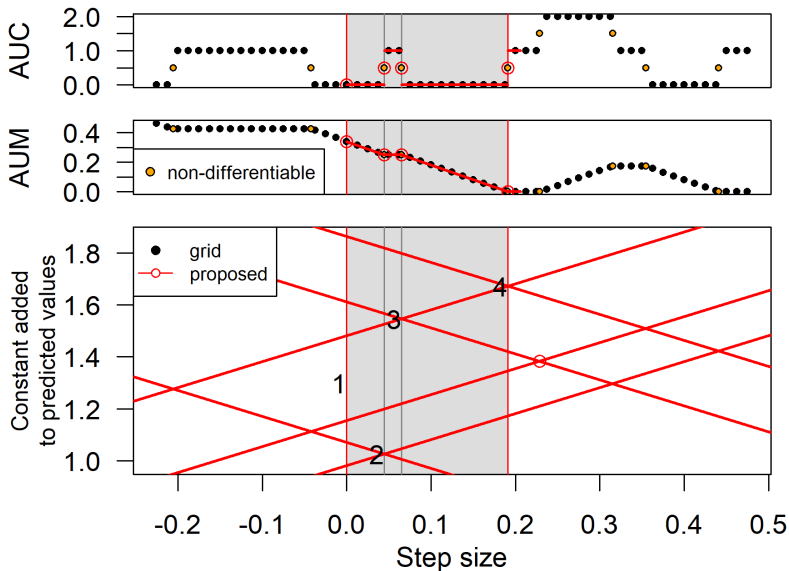
Proposed complete AUC/AUM line search, iteration 3

AUC/AUM values completely known within shaded grey region.



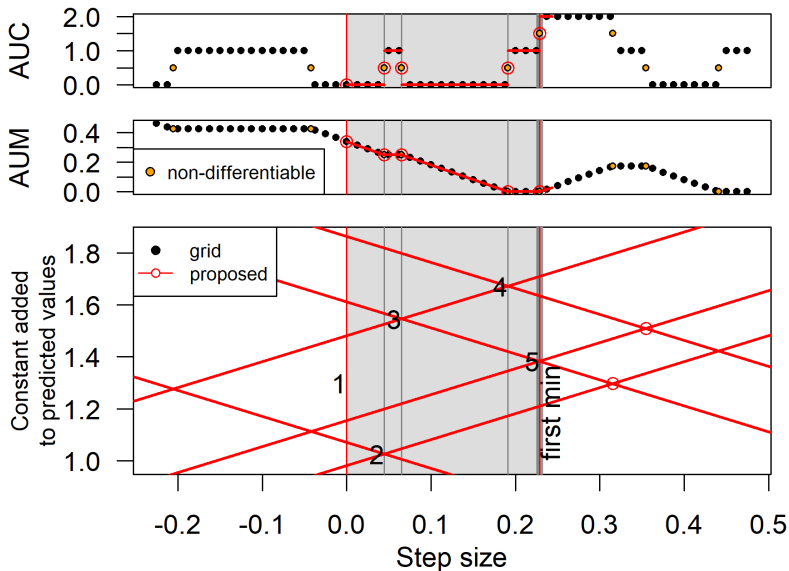
Proposed complete AUC/AUM line search, iteration 4

AUC/AUM values completely known within shaded grey region.



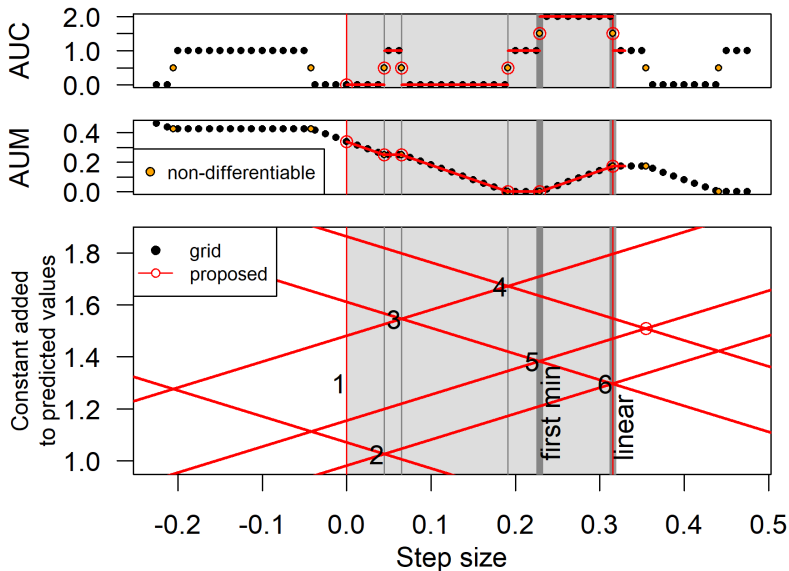
Proposed complete AUC/AUM line search, iteration 5

AUC/AUM values completely known within shaded grey region.



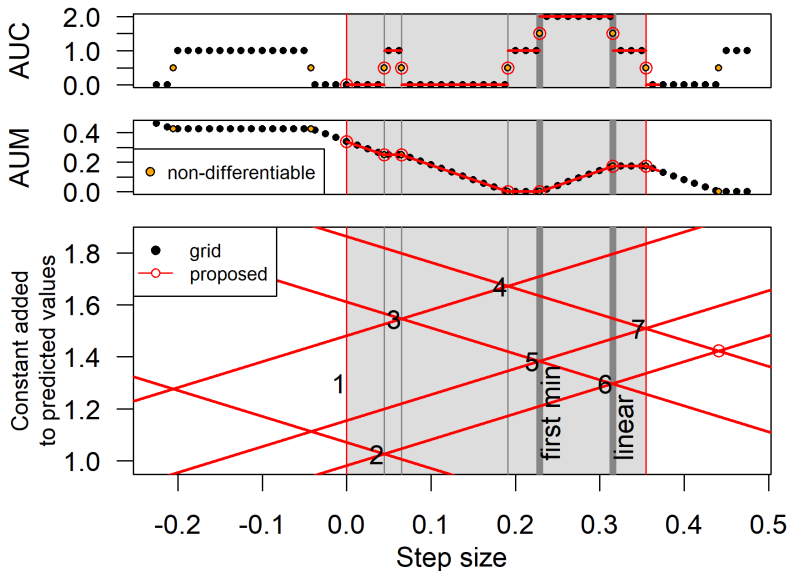
Proposed complete AUC/AUM line search, iteration 6

AUC/AUM values completely known within shaded grey region.



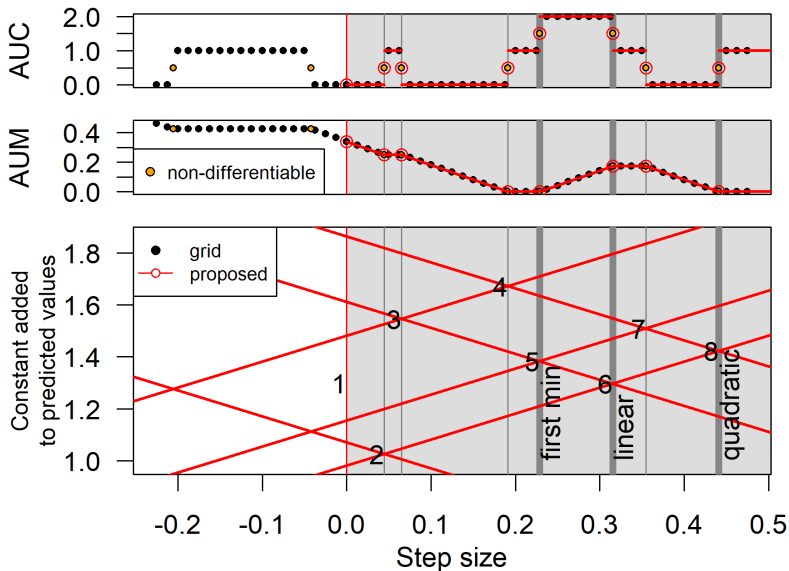
Proposed complete AUC/AUM line search, iteration 7

AUC/AUM values completely known within shaded grey region.



Proposed complete AUC/AUM line search, iteration 8

AUC/AUM values completely known within shaded grey region.



Complexity analysis of proposed algorithm

For N labeled observations, input N threshold line slope/intercept values. Possible next intersection points stored in a C++ STL map (red-black tree, sorted by step size), $O(\log N)$ time insertion, $O(1)$ lookup of next intersection. Worst case $O(N)$ space.

grid: standard grid search. $O(GN \log N)$ time per step, for G grid points.

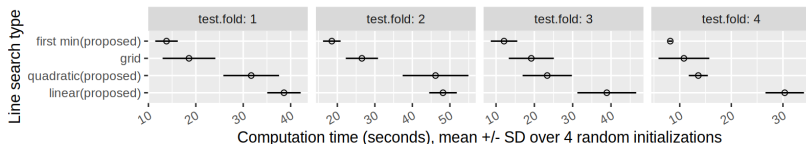
linear(proposed): only first N intersections. $O(N \log N)$ time per step, relatively small step sizes chosen, relatively large number of steps overall in gradient descent.

quadratic(proposed): all $O(N^2)$ intersections. $O(N^2 \log N)$ time per step, large step sizes, small number of steps.

first min(proposed): keep iterating until first AUM increase. Same as quadratic in worst case, but may be faster on average (it was faster than both quadratic and linear for the example on the previous slide).

Proposed search consistently faster than grid search

Analyzed supervised genomic change-point detection data set H3K4me3_TDH_immune ($N = 1073$ to 1248) from UCI Machine Learning Repository, <https://archive.ics.uci.edu/ml/datasets/chipseq>, Train/test splits defined via 4-fold CV, linear model initialized by minimizing regularized convex loss (surrogate for label error, Hocking *et al.* ICML 2013), keep doing AUM rate gradient descent steps (with line search) until subtrain loss stops decreasing.



first min(proposed): keep iterating until first AUM increase.

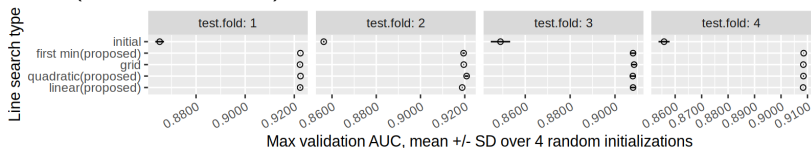
grid: search over step size $\in \{10^{-9}, 10^{-8}, \dots, 10^1, 10^0\}$.

quadratic(proposed): all line search iterations.

linear(proposed): only first N line search iterations.

Proposed search has similar accuracy as grid search

Analyzed supervised genomic change-point detection data set H3K4me3_TDH_immune ($N = 1073$ to 1248) from UCI Machine Learning Repository, <https://archive.ics.uci.edu/ml/datasets/chipseq>, Train/test splits defined via 4-fold CV, linear model initialized by minimizing regularized convex loss (surrogate for label error, Hocking *et al.* ICML 2013), keep doing AUM rate gradient descent steps (with line search) until subtrain loss stops decreasing.



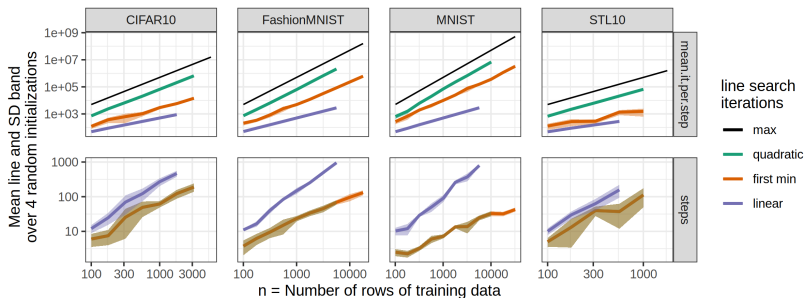
first min(proposed): keep iterating until first AUM increase.

grid: search over step size $\in \{10^{-9}, 10^{-8}, \dots, 10^1, 10^0\}$.

quadratic(proposed): all line search iterations.

linear(proposed): only first N line search iterations.

Asymptotic time complexity analysis



max: theoretical worst case, $N(N - 1)/2$ iterations.

quadratic(proposed): all line search iterations.

first min(proposed): keep iterating until first AUM increase (same number of steps/solution as quadratic, but asymptotically faster/smaller slope).

linear(proposed): only first N line search iterations.

Problem Setting 1: ROC curves for evaluating supervised binary classification algorithms

Proposed AUM loss, a differentiable surrogate for ROC curve optimization (JMLR'23)

Problem setting 2: ROC curves for evaluating supervised changepoint algorithms

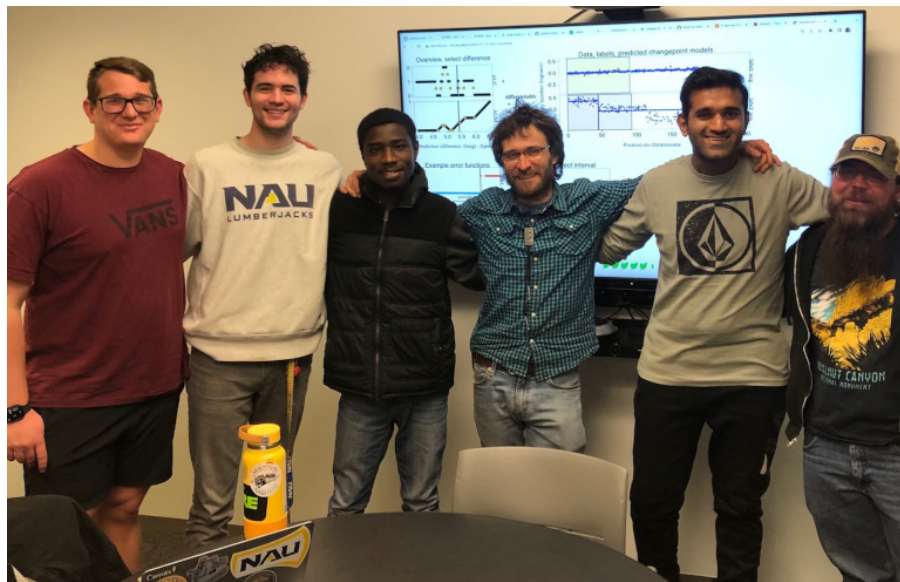
Proposed complete line search algorithm for optimizing AUM and AUC (in progress)

Discussion and Conclusions

Discussion and Conclusions

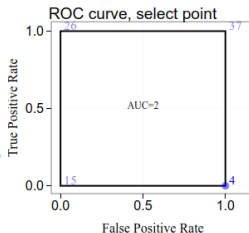
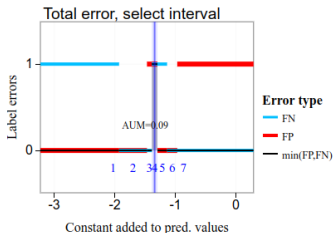
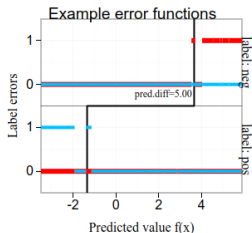
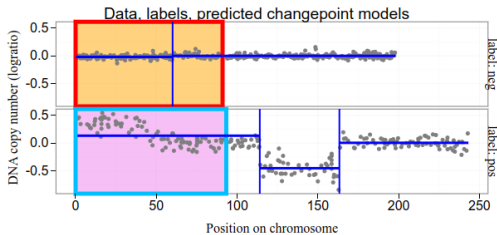
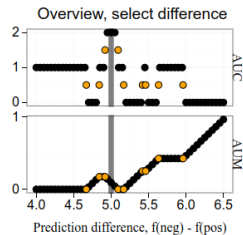
- ▶ Area Under the ROC Curve (AUC) is used to evaluate binary classification and changepoint detection algorithms.
- ▶ Hocking, Hillman, *Journal of Machine Learning Research* (2023), proposed $\text{AUM} = \text{Area Under Min}(\text{FP}, \text{FN})$, a new differentiable surrogate loss for AUC optimization.
- ▶ In this talk we proposed new gradient descent algorithms with efficient complete line search, for optimizing AUM/AUC.
- ▶ Empirical results provide evidence that proposed complete line search is consistently faster than grid search, and has comparable accuracy (in terms of max validation AUC).
- ▶ Line search implemented in R/C++:
<https://cloud.r-project.org/web/packages/aum/> (R/C++ line search)
- ▶ Future work: line search for non-linear prediction functions such as neural networks with ReLU activation?

Thanks to co-author Jadon Fowler! (second from left)



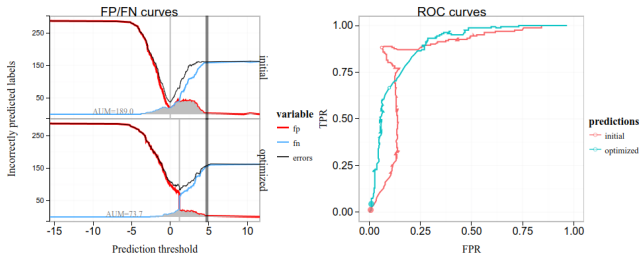
Contact: toby.hocking@nau.edu

AUC can be greater than one for change-point problems with non-monotonic error functions

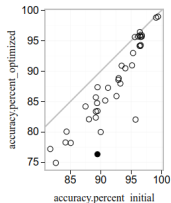


<https://nhintruong.github.io/figure-aum-convexity-interactive>

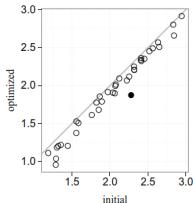
Initial/optimized AUC/AUM for change-point problems



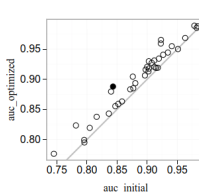
Percent correctly predicted labels



Log[Area under Min(FP,FN) + 1]



Area under the ROC curve



<https://tdhock.github.io/2023-11-21-auc-improved>