Introduction
ooooo

Background
ooooooo

Methodology
oooooo

Results
oooooooooooo

References

# Sub-Quadratic AUC Optimization

Kyle Rust

Northern Arizona University

December 9, 2022

Introduction
ooooo

Background
ooooooo

Methodology
oooooo

Results
ooooooooooooo

References

Introduction

Background

Methodology

Results

**Introduction**
○●○○○

Background
○○○○○○○○

Methodology
○○○○○○

Results
○○○○○○○○○○○○○

References

# Introduction

- ▶ Binary classification aims to learn some function $f(x)$ to predict a positive or negative label
- ▶ To learn $f(x)$ usually involves some sort of optimization of an objective function
- ▶ If the the goal is to minimize the objective function it is referred to as a loss function
- ▶ The focus of this presentation will be on two specific loss functions, the square and squared hinge loss functions

Kyle Rust

Northern Arizona University

Sub-Quadratic AUC Optimization

# Binary Classification Examples

- ▶ Binary classification problems appear in all sorts of different domains
- ▶ *Ex:*
  - ▶ Classifying emails as spam or not spam
  - ▶ Determine whether a image contains a cat or a dog
  - ▶ Identifying invasive lake trout in Yellowstone Lake from LiDAR data
- ▶ Being able to compute the squared-hinge loss in sub-quadratic time would be very valuable
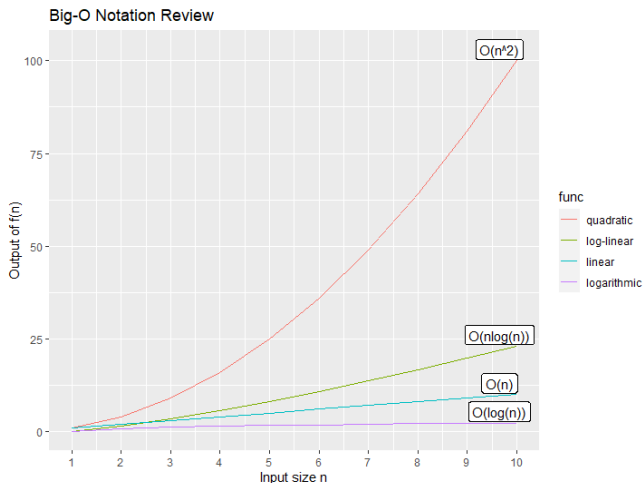- ▶ Particularly for solving problems that have high class imbalance

Kyle Rust                                           Northern Arizona University

Sub-Quadratic AUC Optimization

Introduction
ooooo
Background
ooooooo
Methodology
oooooo
Results
ooooooooooooo
References

# Reason For Research Study

**Naïve Square Loss**

1: Input: Predictions $\hat{y}_1, \ldots, \hat{y}_n \in \mathbb{R}$, labels
   $y_1, \ldots, y_n \in \{-1, 1\}$, margin size $m \geq 0$.
2: Initialize loss to zero
3: $\mathcal{I}^+ = \{i \mid y_i = 1\}$
4: $\mathcal{I}^- = \{i \mid y_i = -1\}$
5: **for** $j \in \mathcal{I}^+$ **do**:
6:     **for** $k \in \mathcal{I}^-$ **do**:
7:        $z \leftarrow \hat{y}_j - \hat{y}_k$
8:        $loss \mathrel{+}= (m - z)^2$
9: Output: total loss: $loss$

This double for loop results in a time complexity of $O(n^2)$

Kyle Rust          Northern Arizona University

Sub-Quadratic AUC Optimization

# Big-O Notation



Big-O Notation Review

Kyle Rust                                          Northern Arizona University
Sub-Quadratic AUC Optimization

# Research Question

How beneficial is it to compute the square and squared
hinge loss in sub-quadratic time?

Introduction
○○○○○

Background
●○○○○○○○

Methodology
○○○○○○

Results
○○○○○○○○○○○○○

References

# Objective Function

- ▶ Machine learning, particularly deep-learning, involves one form of optimization

- ▶ The goal is to either maximize or minimize some function $f(x)$ by altering its input $x$

- ▶ This function $f(x)$ is referred to as an *objective function*

- ▶ If the goal is to minimize $f(x)$ this function is specifically referred to as a *loss function*

Kyle Rust
Northern Arizona University

Sub-Quadratic AUC Optimization

Introduction
ooooo

Background
o●oooooo

Methodology
oooooo

Results
oooooooooooooo

References

# Gradient Descent

- ▶ One way to intelligently select $x \rightarrow f(x)$ is to use derivative information
- ▶ The derivative with respect to the loss is taken and that information is used to traverse to the minimum
- ▶ The derivative points in the direction of steepest *ascent*, so we go in the opposite direction for steepest *descent*
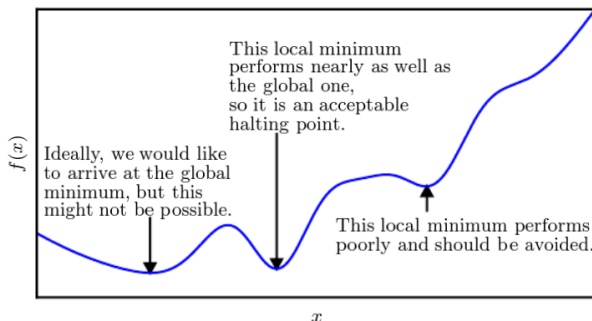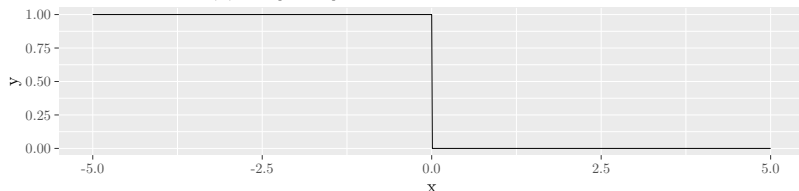


*Figure: (Goodfellow et al., 2016)*

Introduction
00000

Background
00●00000

Methodology
000000

Results
0000000000000

References

# Convex Surrogate

▶ In order to make use of gradient for optimization a function must be differentiable and the derivative must be non-zero

▶ A function that can be optimized in place of the objective function is called a convex surrogate

▶ A function is convex if its local minima and maxima are also global minima and maxima

Zero-One Loss $\ell(z) = I[z < 0]$

Kyle Rust                                                                 Northern Arizona University

Sub-Quadratic AUC Optimization

Introduction
ooooo

**Background**
oooo●oooo

Methodology
oooooo

Results
oooooooooooo

References

# Cross Validation

▶ We want to develop a machine learning model that can make predictions on new data

▶ To do this we have to assume our new data is similar to the data we train on
  1. In statistics this is called independent and identically distributed

▶ One method for doing so is to use the K-Fold methodology

▶ To do this the data is first separated in train and test sets
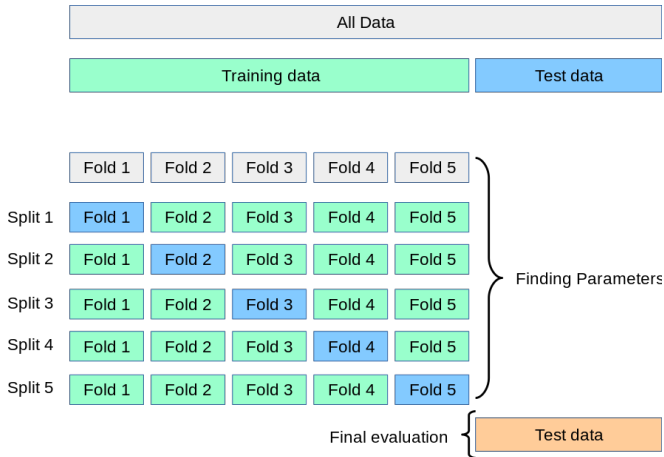
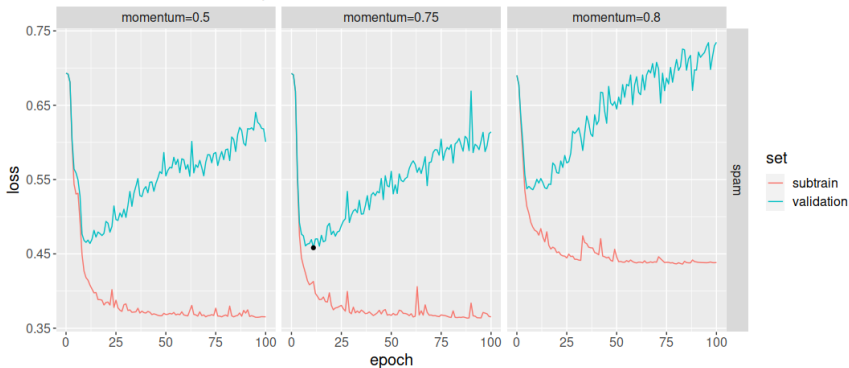▶ The train set is then split again into the subtrain and validation set
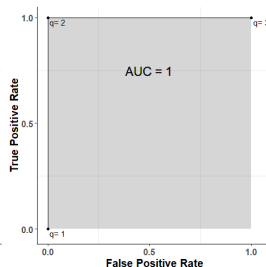
# K-Fold



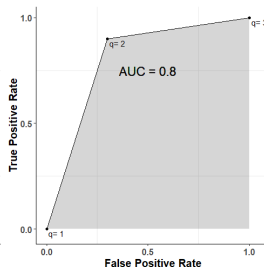Figure: (Pedregosa et al., 2011)

Kyle Rust                                                    Northern Arizona University

Sub-Quadratic AUC Optimization

# Regularization



Loss As A Function of Epochs For Different Momentum Values

Kyle Rust                                        Northern Arizona University
Sub-Quadratic AUC Optimization

Introduction
ooooo

**Background**
oooooo●o

Methodology
oooooo

Results
ooooooooooooo

References

# Receiver Operating Characteristic Curve

▶ The Receiver Operating Characteristic Curve (ROC) is a plot of False Positive Rate (FPR) vs True Positive Rate (TPR)

Kyle Rust

Northern Arizona University

Sub-Quadratic AUC Optimization

Introduction
ooooo
**Background**
ooooooo●
Methodology
oooooo
Results
ooooooooooooo
References

# Context Within Greater Discipline

▶ Benefits to maximizing AUC:
  1. Optimizing for a common machine learning metric can lead to higher model performance
  2. ROC AUC is more adept at handling highly unbalanced data than other common metrics

▶ Challenges:
  1. Finding a surrogate function for ROC AUC is difficult
  2. Small changes to the learning model, can result in large changes to the ROC AUC

Introduction
ooooo

Background
ooooooo

Methodology
●ooooo

Results
oooooooooooo

References

## Related Work Summary

| Paper | Degree | Hinge | Proof | Solution |
|---|---|---|---|---|
| Pahikkala et al. | Square | False | False | Functional |
| Joachims | Linear | True | True | Functional |
| Calders and Jaroszewicz | Polynomial | False | True | Functional |
| Ying et al. | Square | False | True | Min-Max |
| Yuan et al. | Square | True | True | Min-Max |
| This Work | Square | True | True | Functional |

Maximizing AUC: (Bamber, 1975), (Herschtal and Raskutti, 2004)

Re-weighting/Sorting Algorithms: (Calders and Jaroszewicz, 2007),(Ying et al., 2016), (Yuan et al., 2020)

Pairwise Algorithms: (Pahikkala et al., 2009), (Joachims, 2005)

Introduction
ooooo
Background
ooooooo
Methodology
o●oooo
Results
ooooooooooooo
References

# Loss Functions That Sum Over Examples

▶ To solve a binary classification, we want to learn some function $f : \mathbb{R}^p \to \mathbb{R}$, where $p$ is the number of features

▶ The real-valued predictions are computed by $\hat{y}_i = f(x_i)$

▶ For every observation $\mathcal{L}(f) = \sum_{i=1}^n \ell[y_i f(\mathbf{x}_i)]$

▶ Values where $y_i f(x_i) > 0$ result in correct predictions in accordance with the labels, $y_i f(x_i) < 0$ results in incorrect predictions.

Kyle Rust                                                                Northern Arizona University

Sub-Quadratic AUC Optimization

Introduction
ooooo

Background
ooooooo

Methodology
ooo●oo

Results
ooooooooooooo

References

# Pairwise Loss Functions

▶ Pairwise loss functions sum over all pairs of the positive examples and negative examples

▶ $\mathcal{L}(f) = \sum_{j \in \mathcal{I}^+} \sum_{k \in \mathcal{I}^-} \ell[f(\mathbf{x}_j) - f(\mathbf{x}_k)]$

▶ Positive pairwise distances $f(\mathbf{x}_j) - f(\mathbf{x}_k) > 0$ result in correctly ranked pairs, while negative pairwise distances $f(\mathbf{x}_j) - f(\mathbf{x}_k) < 0$ result in incorrectly ranked pairs

Introduction
ooooo

Background
ooooooo

Methodology
ooo●oo

Results
ooooooooooooo

References

## Functional Square Loss

Basic definition of the square loss:

1. $\ell(z) = (m - z)^2 \; z = \hat{y}_j - \hat{y}_k \; m = \text{margin parameter}$

Substitute the definition of the square loss into the definition of the pairwise loss:

2. $\sum_{j \in \mathcal{I}^+} \ell(\hat{y}_j - \hat{y}_k) = \sum_{j \in \mathcal{I}^+} (m - \hat{y}_j + \hat{y}_k)^2$

Group margin and predicted values, expand the terms:

3. $= \sum_{j \in \mathcal{I}^+} ((m - \hat{y}_j) + \hat{y}_k)((m - \hat{y}_j) + \hat{y}_k)$
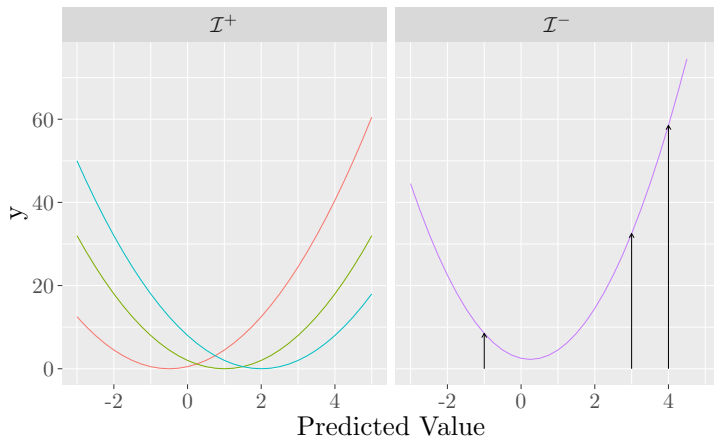
Factor the polynomial:

4. $= \sum_{j \in \mathcal{I}^+} \underbrace{1}_{a} \hat{y}_k^2 + \underbrace{2(m - \hat{y}_j)}_{b} \hat{y}_k + \underbrace{(m - \hat{y}_j)^2}_{c}$

Substitute to coefficients:
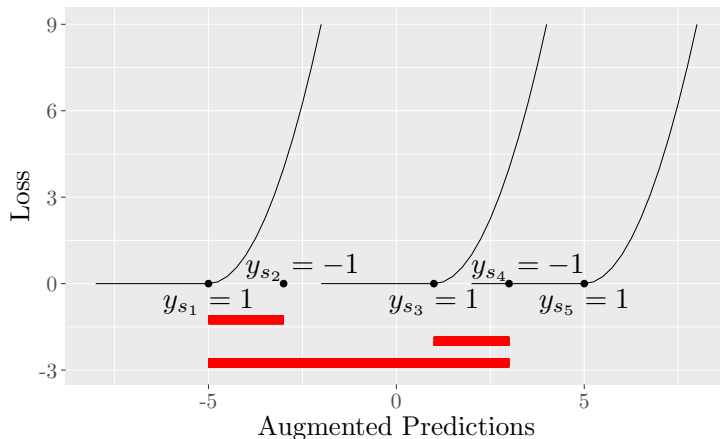
5. $= a^+ \hat{y}_k^2 + b^+ \hat{y}_k + c^+$

Introduction
00000

Background
00000000

Methodology
00000●0

Results
0000000000000

References

# Geometric View of Square Loss



Square Loss Accrues Loss

Kyle Rust

Northern Arizona University

# Geometric View of Functional Squared Hinge



Sorted Predictions Accrues Loss

Introduction
ooooo
Background
ooooooo
Methodology
oooooo
Results
●oooooooooooo
References

## Monsoon Cluster

- ▶ Timing experiments were carried out on a personal machine with a Intel(R) Core(TM) i5-8600 CPU @ 3.10GHz CPU

- ▶ All model training experiments were carried out on NAU's computing cluster Monsoon

- ▶ Monsoon consists of 4076 cores, 26TB of memory, and 27 NVIDIA GPUs (HPC, 2021)

- ▶ These model performance experiments were computed using AMD EPYC 7542 CPUs

Introduction
ooooo

Background
ooooooo

Methodology
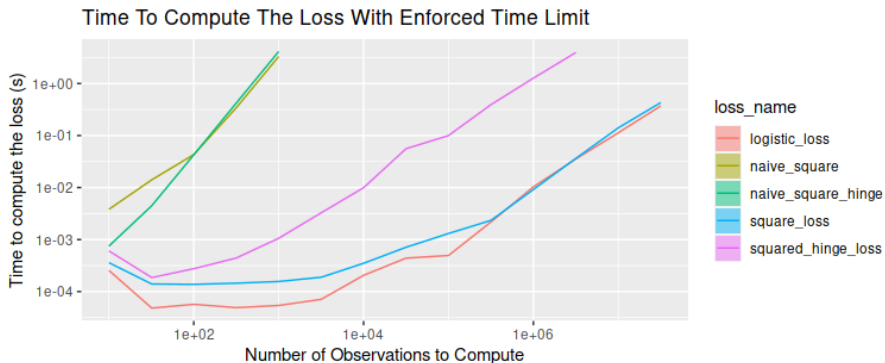ooooo

Results
oooooooooooo

References

# PyTorch

- ▶ PyTorch is an open-source machine learning library
- ▶ PyTorch provides two main pieces of functionality:
    1. Good performance on GPU's
    2. Automatic gradient calculation
- ▶ This is done using a tensor data structure, where derivative information is stored in a graph-like structure
- ▶ It can be implemented by simply calling **backward** on the loss value

# Novel Contributions

1. Improved Time Complexity

2. Increased Possible Batch Size

3. Increased Model Performance

Kyle Rust                                                                Northern Arizona University

Sub-Quadratic AUC Optimization

# Increased Possible Batch Size



Time To Compute The Loss With Enforced Time Limit

Introduction
ooooo

Background
ooooooo

Methodology
ooooo

Results
oooooooooooooo

References

# Data Sets

▶ We trained a model for three data sets
 1. Cat&Dog (Elson et al., 2007)
 2. CIFAR10 (Krizhevsky, 2009)
 3. STL10 (Coates et al., 2011)

▶ The STL10 and CIFAR10 data sets were converted to binary classification data sets

▶ The model was trained on each data set for three ratios of positive to negative examples: 0.1, 0.01, 0.001

▶ To achieve the desired class imbalanced from these sets, positive examples were removed from the train set

▶ The test set was left balanced between positive and negative examples

Kyle Rust
Northern Arizona University

Sub-Quadratic AUC Optimization

Introduction
ooooo

Background
oooooooo

Methodology
oooooo

**Results**
oooooo●oooooo

References

# Splits

▶ Each of the following splits was completed 5 times, with 5 different random initializations

▶ Each data set was first split into 80% and 20%, for training and hold-out test set

▶ From there the train set was split again 80/20, for gradient descent and hyper-parameter selection

Kyle Rust

Northern Arizona University

Sub-Quadratic AUC Optimization

Introduction
00000

Background
00000000

Methodology
000000

Results
000000●000000

References

# Hyper-Parameters

▶ The model was tuned for two different hyper-parameters

1. Batch Size - How many examples are processed in gradient descent
2. Learning Rate - How big of a step is taken for each batch in gradient descent

▶ The batch size was selected from 10, 50, 100, 500, 1000, 5000

▶ The AUCM and logistic losses tested learning rates from $10^{-4}, \dots, 10^2$

▶ The proposed functional squared hinge loss searched over $10^{-4}, \dots, 10^{-1}$

Kyle Rust

Northern Arizona University

Sub-Quadratic AUC Optimization

Introduction
○○○○○

Background
○○○○○○○

Methodology
○○○○○

Results
○○○○○○○●○○○○○

References

# Loss Functions

- ▶ The proposed functional squared hinge loss was compared to a state-of-the-art and industry standard loss functions

- ▶ The AUCM loss proposed by Yuan et al. (2020) serves as the state-of-the-art comparison

- ▶ The logistic (binary cross entropy) loss serves as the basic comparison

- ▶ For the logistic loss, equal weights were used. No special consideration is given to AUC or class imbalances

Introduction
ooooo

Background
ooooooo

Methodology
oooooo

Results
ooooooooo●oooo

References

# Model Architecture

- ▶ The model that was used for training was the ResNet20 architecture (He et al., 2015)
- ▶ The model was initially created to solve the CIFAR10 data set
- ▶ The 20 layers consist of convolutional layers, activation layers, a softmax layer, a fully connected layer, and a reshape
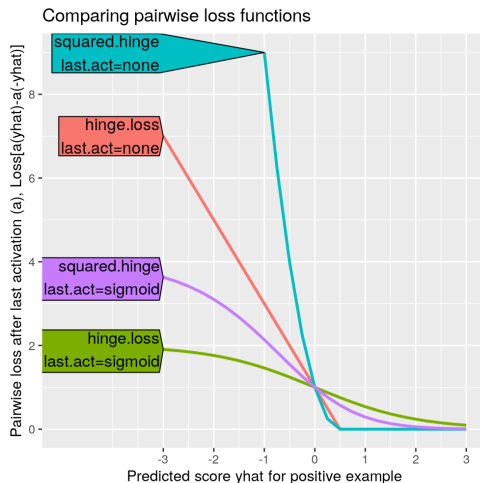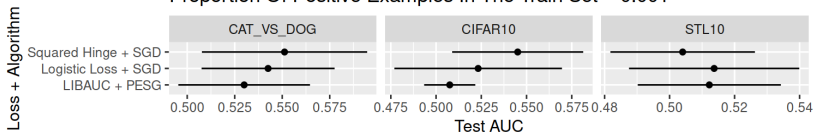- ▶ This resulted in 270,000 parameters to learn

# Last Layer Comparison



*Figure: (Hocking, 2022)*

# Increased Possible Batch Size

| Imratio | Loss Function | CIFAR10 | | STL10 | | Cat&Dog | |
|---------|---------------|---------|---------------|-------|---------------|---------|---------------|
| | | Batch | Learning Rate | Batch | Learning Rate | Batch | Learning Rate |
| 0.1 | Our Square Hinge | 10 | 0.0316 | 10 | 0.0100 | 50 | 0.1000 |
| | LIBAUC | 50 | 0.1000 | 50 | 0.1000 | 50 | 0.1000 |
| | Logistic Loss | 10 | 0.1000 | 50 | 0.1000 | 50 | 1.0000 |
| 0.01 | Our Square Hinge | 10 | 0.0032 | 100 | 0.1000 | 50 | 0.0316 |
| | LIBAUC | 50 | 0.1000 | 1000 | 0.1000 | 100 | 0.1000 |
| | Logistic Loss | 10 | 0.1000 | 1000 | 0.1000 | 100 | 1.0000 |
| 0.001 | Our Square Hinge | **500** | 0.0316 | 10 | 0.0001 | **1000** | 0.3162 |
| | LIBAUC | 100 | 10.0000 | 10 | 0.0001 | 500 | 10.0000 |
| | Logistic Loss | 100 | 1.0000 | 100 | 0.0001 | 100 | 1.0000 |

# Increased Model Performance

Introduction
00000

Background
00000000

Methodology
000000

Results
000000000000●

References

# Conclusion

- ▶ In this presentation algorithms for computing the square and squared hinge loss in $O(n)$ and $O(n \log n)$
- ▶ Background information and related work were discussed for context around the problem
- ▶ A preliminary experiment showing the asymptotic time complexity of the proposed algorithms
- ▶ Demonstrated how it can be advantageous to select larger batch sizes that were not previously feasible
- ▶ Finally it was shown that these new methods can out perform baseline and state-of-the-art loss functions

# Bibliography I

Bamber, D. (1975). The area above the ordinal dominance graph and the area below the receiver operating characteristic graph. *Journal of mathematical psychology*, 12(4):387–415.

Calders, T. and Jaroszewicz, S. (2007). Efficient AUC optimization for classification. In *European Conference on Principles of Data Mining and Knowledge Discovery*, pages 42–53. Springer.

Coates, A., Ng, A., and Lee, H. (2011). An Analysis of Single Layer Networks in Unsupervised Feature Learning. In *AISTATS*. https://cs.stanford.edu/~acoates/papers/coatesleeng_aistats_2011.pdf.

Elson, J., Douceur, J. J., Howell, J., and Saul, J. (2007). Asirra: A captcha that exploits interest-aligned manual image categorization. In *Proceedings of 14th ACM Conference on Computer and Communications Security (CCS)*. Association for Computing Machinery, Inc.

Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep Learning*. Adaptive Computation and Machine Learning series. MIT Press.

He, K., Zhang, X., Ren, S., and Sun, J. (2015). Deep residual learning for image recognition. *CoRR*, abs/1512.03385.

Herschtal, A. and Raskutti, B. (2004). Optimising area under the ROC curve using gradient descent. In *Proceedings of the twenty-first international conference on Machine learning*, page 49.

Hocking, T. D. (2022). Max generalized auc. https://github.com/tdhock/max-generalized-auc.

HPC, N. H. P. C. (2021). Monsoon details. Northern Arizona University.

Joachims, T. (2005). A support vector method for multivariate performance measures. In *Proceedings of the 22nd international conference on Machine learning*, pages 377–384.

Kyle Rust                                                Northern Arizona University

Sub-Quadratic AUC Optimization

Introduction
ooooo

Background
ooooooo

Methodology
oooooo

Results
ooooooooooooo

References

# Bibliography II

Krizhevsky, A. (2009). Learning multiple layers of features from tiny images. DOI:10.1.1.222.9220.

Pahikkala, T., Tsivtsivadze, E., Airola, A., Järvinen, J., and Boberg, J. (2009). An efficient algorithm for learning to rank from preference graphs - machine learning.

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.

Ying, Y., Wen, L., and Lyu, S. (2016). Stochastic online AUC maximization. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 451–459.

Yuan, Z., Yan, Y., Sonka, M., and Yang, T. (2020). Robust Deep AUC Maximization: A New Surrogate Loss and Empirical Studies on Medical Image Classification. Preprint arXiv:2012.03173.