

# Boosting Methods

Jesse Ellin, Ohad Nir, Ritvik Vaish, Thomas Su





# Stochastic Gradient Boosting

Jerome Friedman - Departments of Statistics and Stanford Linear Accelerator Center

Published Feb 28, 2002 in Computational Statistics and Data Analysis



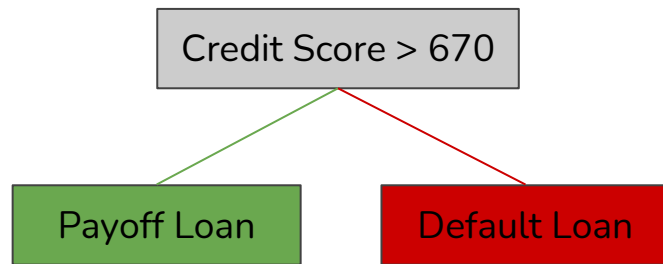
# **What are Boosted Methods?**





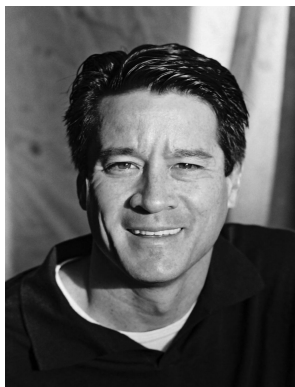
# Weak Learnability and Weak Learners

- Put simply, a predictor that performs slightly better than random guessing.
- Unlike strong learners, we can't guarantee convergence.
- Examples of Weak Learners:
  - Stump
  - Shallow network
- Examples of Strong Learners:
  - ERM using (Stochastic) Gradient Descent
  - Multi-layered Neural Networks



# The History of the Boosting Paradigm

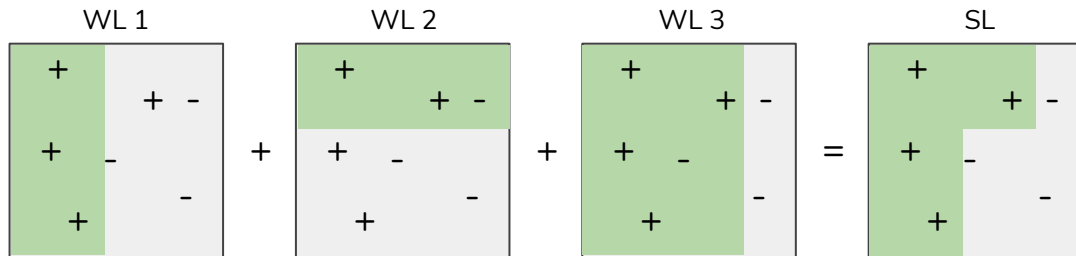
- Paradigms driving question:
  - Can an efficient weak learner can be “boosted” into an efficient strong learner?
  - Question proposed by Michael Kearns and Leslie Valiant in 1988.
    - Paper: Cryptographic Limitations on Learning Boolean Formulae and Finite Automata.
- Solutions:
  - Solved by Robert Schapire in 1990.
    - Paper: The strength of weak learnability.
  - Adaboost, the first practical algorithm for boosting was proposed by Robert Schapire and Yoav Freund in 1995 and later won the Godel Prize.
    - Paper: A Short Introduction to Boosting.





# What is Boosting?

- Trains weak learners to be strong learners.
  - Does this through combining weak learners.
- Addresses two issues:
  - Bias-complexity tradeoff:
    - The error of an ERM could be decomposed into approximation error and estimation error.
      - Approximation error - Determined by the quality of our prior knowledge.
      - Estimation error - Error due to overfitting.
    - The larger the hypothesis class, the smaller the approximation error and the larger the estimation error.
    - Boosted methods allow the learner to control this tradeoff.
  - Computational complexity of learning:
    - In general, complex predictors are better predictors, but they are harder to train.
    - What if we take a linear combination of simple predictors?



# Gradient Boosting vs Stochastic Boosting

## Regular Boosting

- When training, we take the full training data set per epoch.
- More likely to overfit to the training data.
- Long training time.

	Algorithm 1: Gradient_TreeBoost
1	$F_0(\mathbf{x}) = \arg \min_{\gamma} \sum_{i=1}^N \Psi(y_i, \gamma)$
2	For $m = 1$ to $M$ do:
3	$\tilde{y}_{im} = - \left[ \frac{\partial \Psi(y_i, F(\mathbf{x}_i))}{\partial F(\mathbf{x}_i)} \right]_{F(\mathbf{x})=F_{m-1}(\mathbf{x})}, i = 1, N$
4	$\{R_{lm}\}_1^L = L - \text{terminal node } tree(\{\tilde{y}_{im}, \mathbf{x}_i\}_1^N)$
5	$\gamma_{lm} = \arg \min_{\gamma} \sum_{\mathbf{x}_i \in R_{lm}} \Psi(y_i, F_{m-1}(\mathbf{x}_i) + \gamma)$
6	$F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \nu \cdot \gamma_{lm} 1(\mathbf{x} \in R_{lm})$
7	endFor

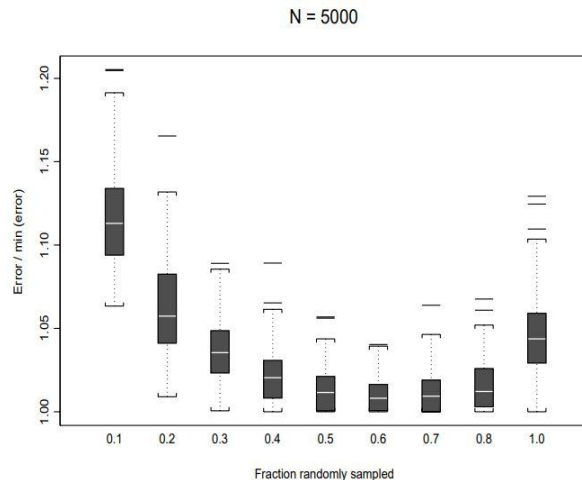
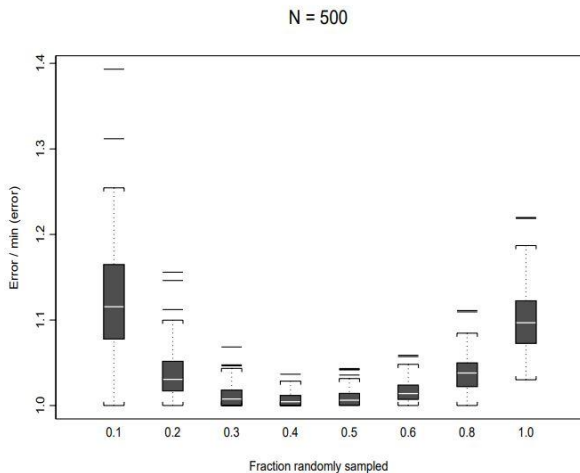
## Stochastic Boosting

- When training model, we sample random subsets of data (without replacement) over multiple epochs.
- Increases variance which make it less likely to overfit.
- Faster training time.

	Algorithm 2: Stochastic Gradient_TreeBoost
1	$F_0(\mathbf{x}) = \arg \min_{\gamma} \sum_{i=1}^N \Psi(y_i, \gamma)$
2	For $m = 1$ to $M$ do:
3	$\{\pi(i)\}_1^N = \text{rand\_perm}\{i\}_1^N$
4	$\tilde{y}_{\pi(i)m} = - \left[ \frac{\partial \Psi(y_{\pi(i)}, F(\mathbf{x}_{\pi(i)}))}{\partial F(\mathbf{x}_{\pi(i)})} \right]_{F(\mathbf{x})=F_{m-1}(\mathbf{x})}, i = 1, \tilde{N}$
5	$\{R_{lm}\}_1^L = L - \text{terminal node } tree(\{\tilde{y}_{\pi(i)m}, \mathbf{x}_{\pi(i)}\}_1^{\tilde{N}})$
6	$\gamma_{lm} = \arg \min_{\gamma} \sum_{\mathbf{x}_{\pi(i)} \in R_{lm}} \Psi(y_{\pi(i)}, F_{m-1}(\mathbf{x}_{\pi(i)}) + \gamma)$
7	$F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \nu \cdot \gamma_{lm} 1(\mathbf{x} \in R_{lm})$
8	endFor

# Papers Results on Regression Problem

- Used M\_TreeBoost procedure (Huber M-regression Loss).
  - Was found to be the best regression procedures in Friedman (1999).
- Tests were done on a small dataset with  $N=500$  data points and a larger dataset with  $N=5000$  data points.
- Different degrees of randomness were applied.
  - Controlled through the fraction  $f=\tilde{N}/N$ .

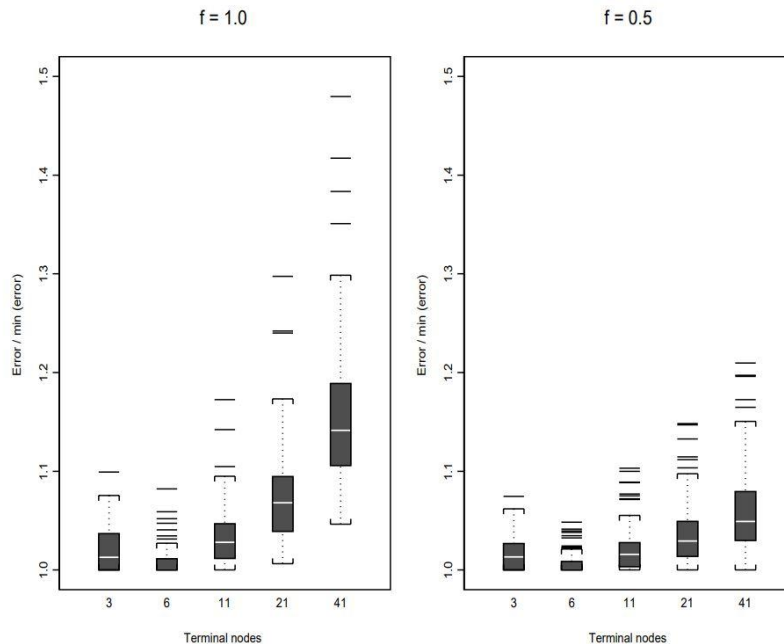






# More Results on Regression Problem

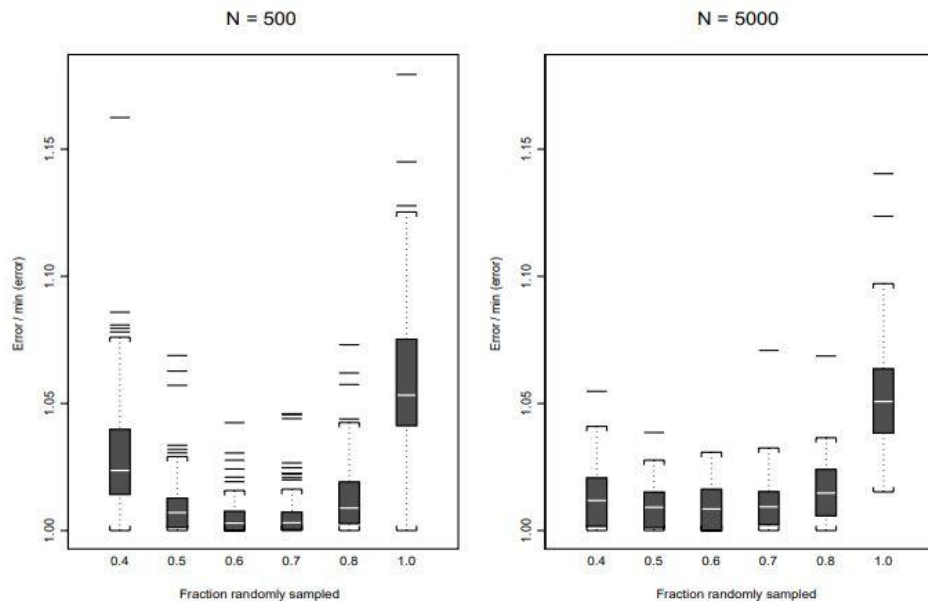
- Used  $N=500$  dataset using M\_TreeBoost procedure.
- Shows Error over different number of terminal nodes (leaf nodes) base learners.





# Papers Results on Classification Problem

- Experiment was done on a dual classification problem.
- Loss criteria was on twice binomial negative log-likelihood.
- Sampled an equal amount of data for both classes.





# The Experiment





# Our Motivation

We want to compare Adaboost, a known and standard Boosting architecture, to a stochastic Boosting architecture to see how adding stochasticity improves or hinders the performance of Boosting.

This will hopefully give us a stronger intuition of how Boosting performs.



# Our Experiment

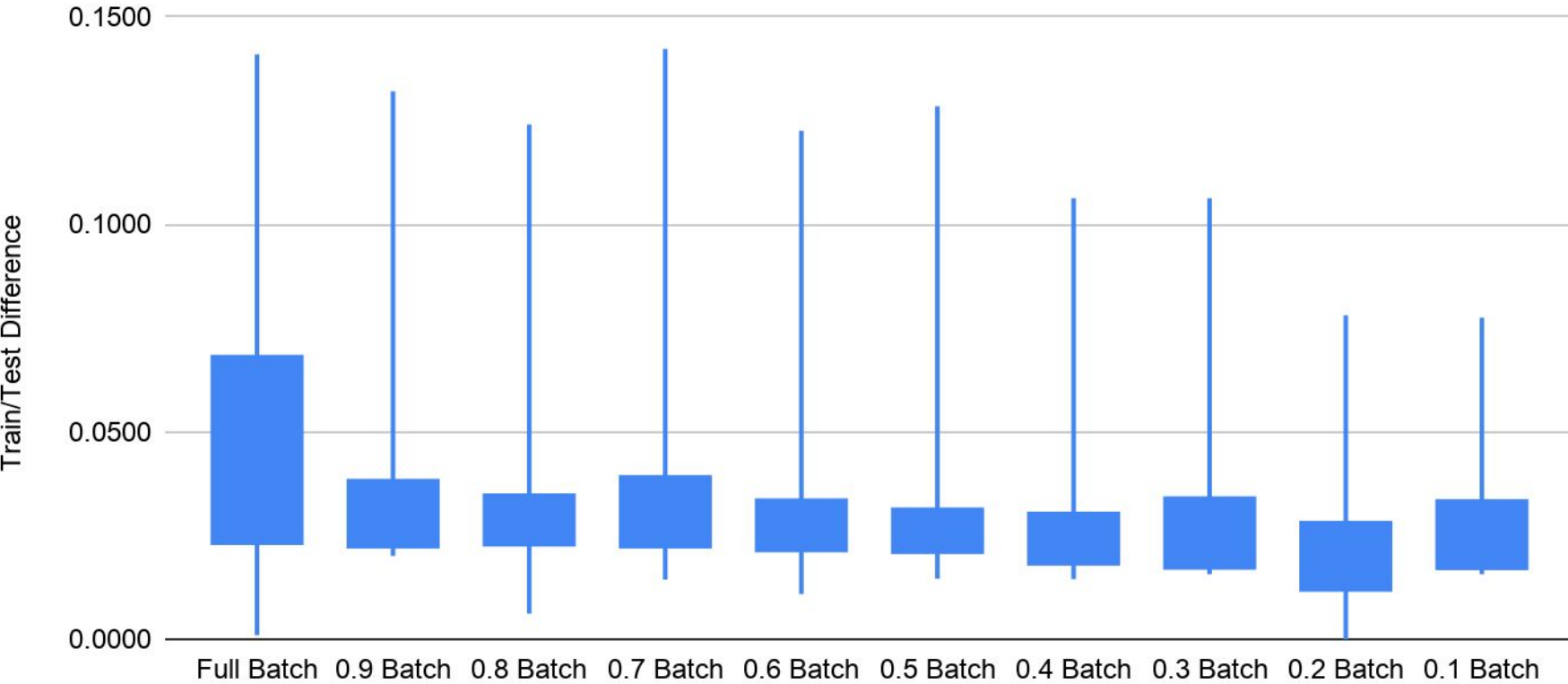
- Experiment with the behavior of boosting with different training data sizes and tree sizes.
  - See if adding stochastic processes results in similar improvement to the paper.
  - See how changing tree depth impacts overfitting trends
  - Is there a general change in runtime (determined by relative training times) when we change batch size or tree size
- Dataset:
  - Dataset 1 - Sign-Language MNIST
  - Dataset 2 - Chinese MNIST
  - Dataset 3 - Fashion MNIST



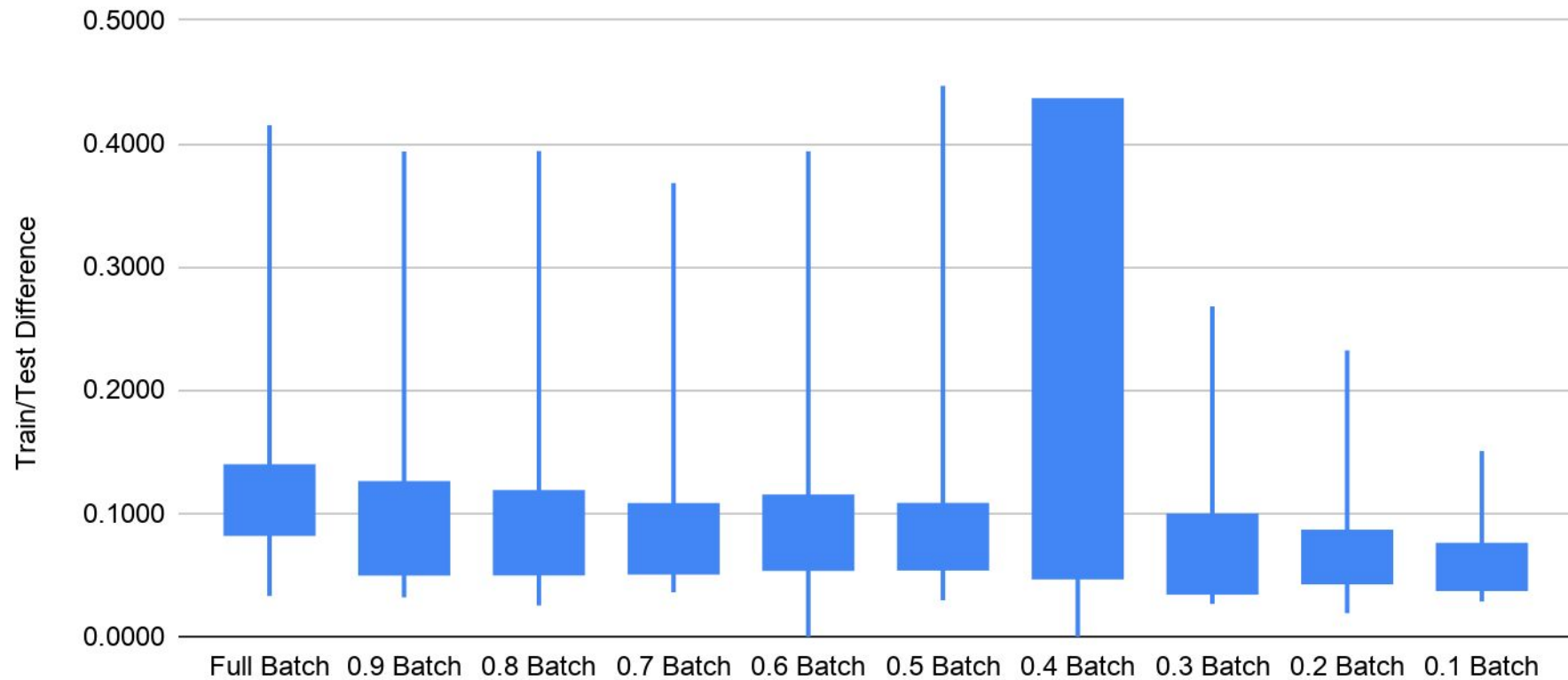
# Results



Dataset 1 - Batch Size Generalization Gap over Different Minibatch Sizes

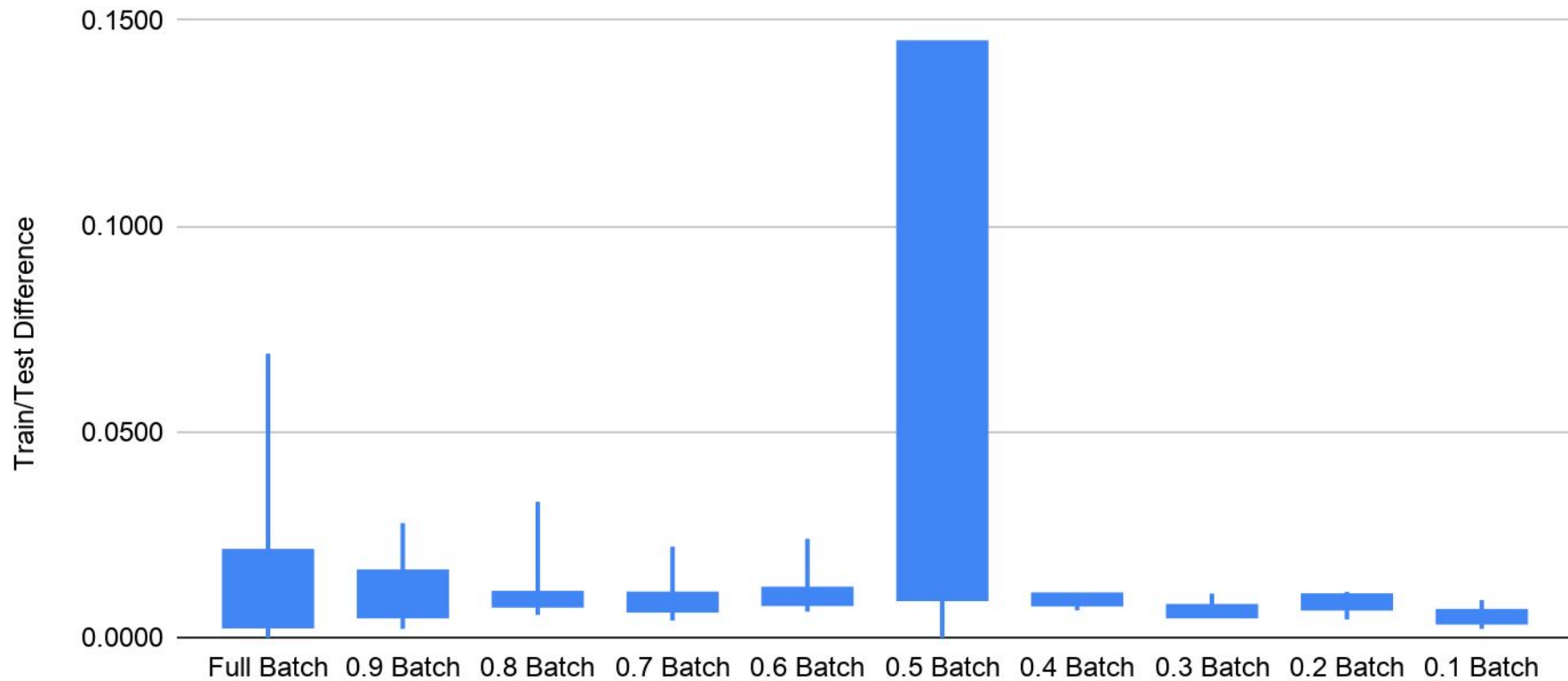


## Dataset 2 - Batch Size Generalization Gap over Different Minibatch Sizes

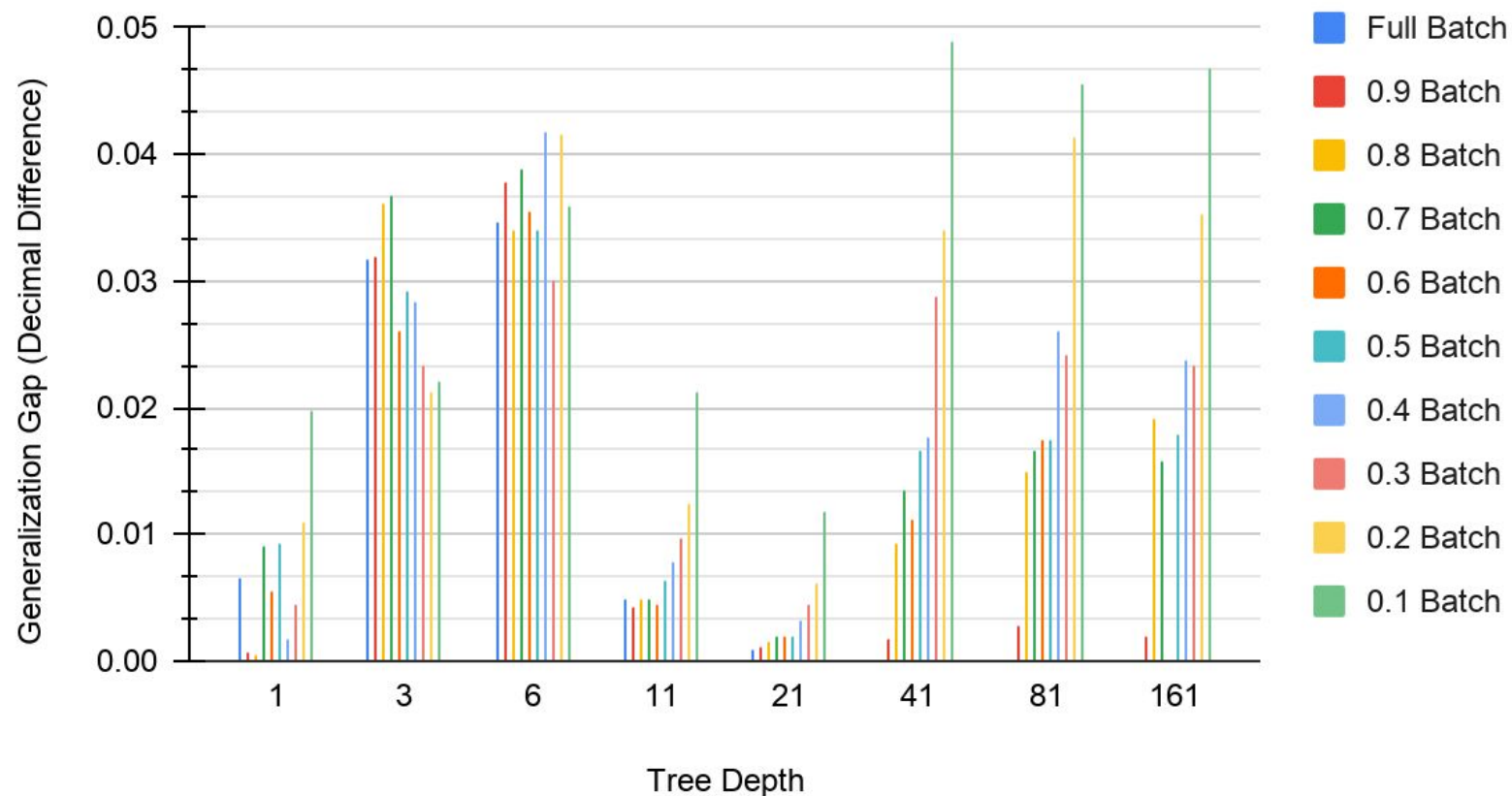




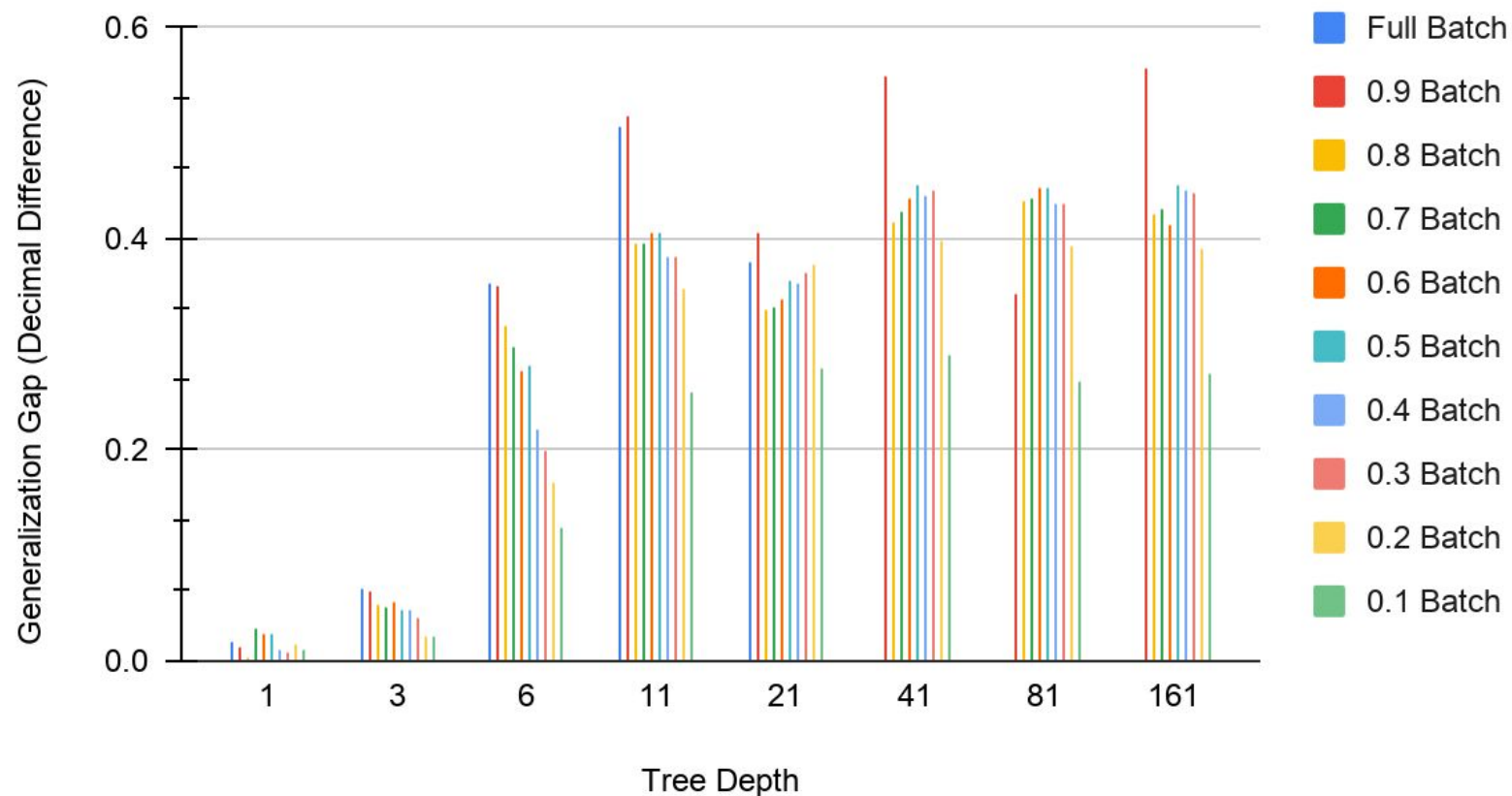
## Dataset 3 - Batch Size Generalization Gap over Different Minibatch Sizes



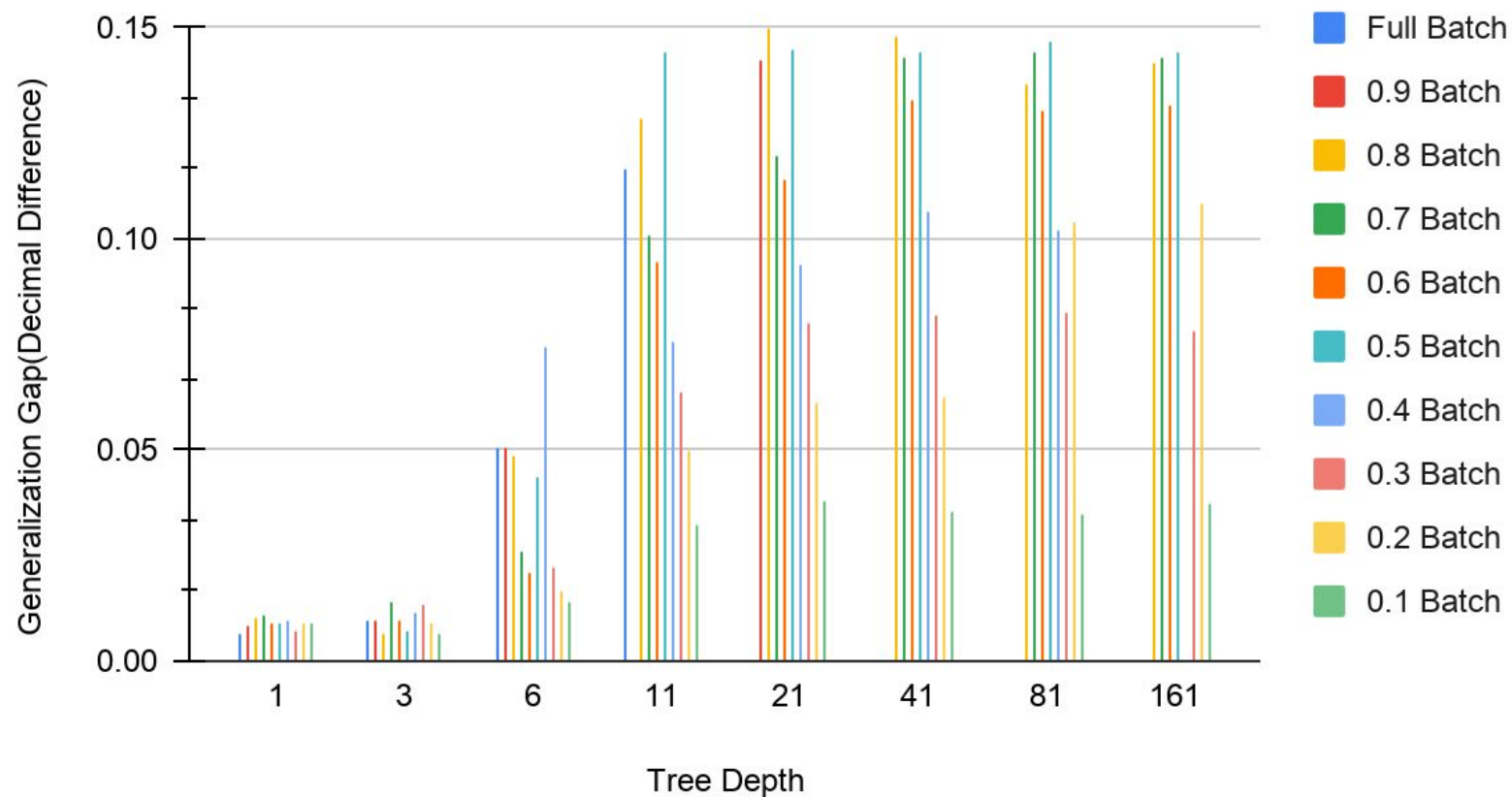
# Dataset 1 Generalization Gap



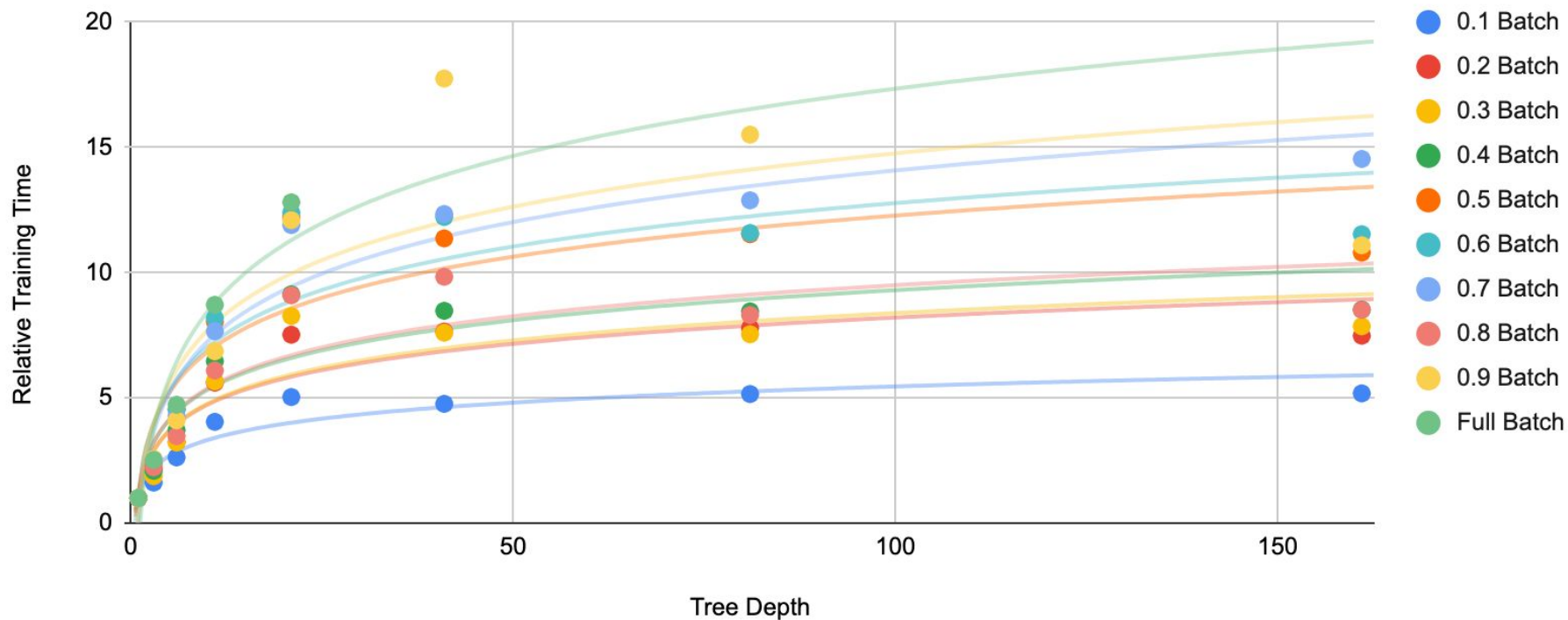
## Dataset 2 Generalization Gap



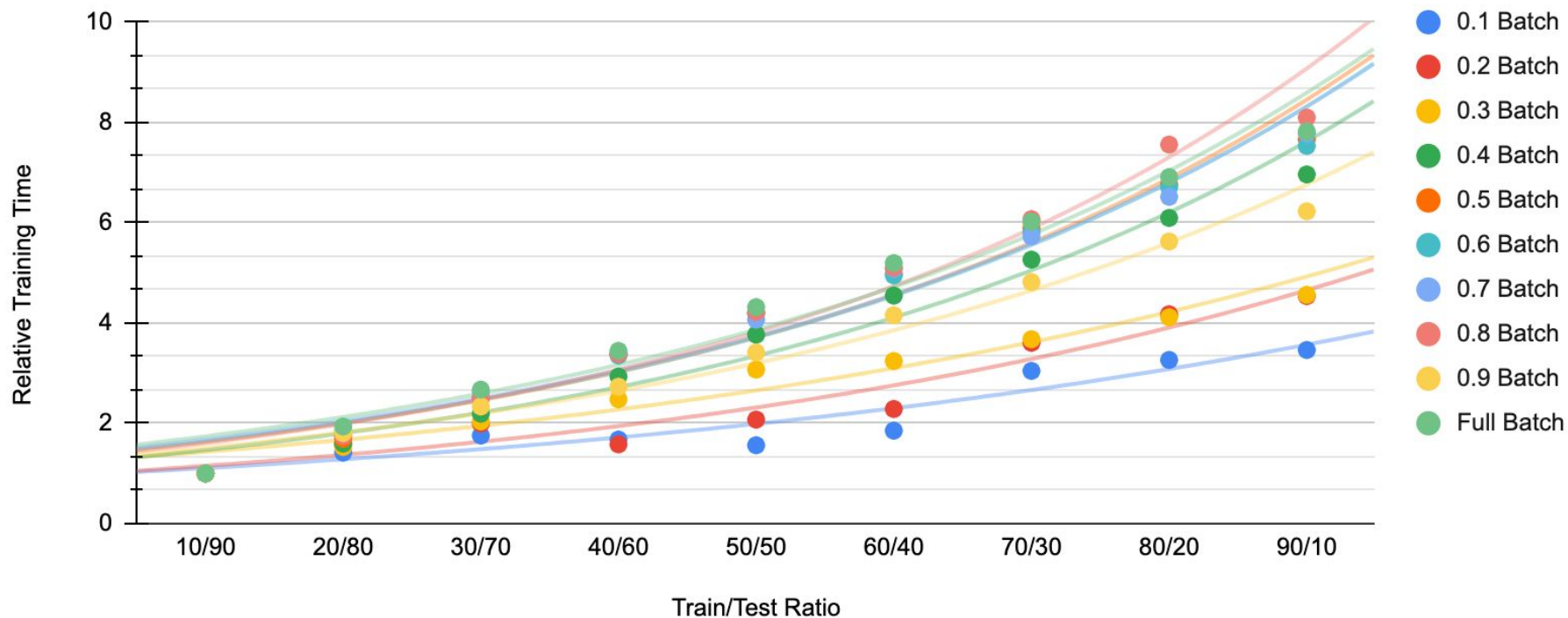
## Dataset 3 Generalization Gap



Dataset 1 - Tree Depth VS Relative Training Time - Log Trendline



Dataset 1 - Train/Test Ratio Vs Relative Training Time - Exponential Trendline





# Homework





# Homework Instruction

- Implementing a framework for a multiclass version of Adaboost.
  - Stagewise Additive Modeling using a Multi-class Exponential loss function (SAMME).
  - Stochastic
- Train your framework on two data sets:
  - Speaker Accent Recognition DataSet.
  - Deepfakes: Medical Image Tamper Detection DataSet.
- Will use two kinds of weak learners.
  - Gaussian Naive Bayes
  - Shallow Decision Tree
- Compare:
  - Regular boosting and stochastic boosting.
  - Weak learners



# SAMME Multiclass AdaBoost Algorithm

## Algorithm 2 SAMME

1. Initialize the observation weights  $w_i = 1/n$ ,  $i = 1, 2, \dots, n$ .

2. For  $m = 1$  to  $M$ :

(a) Fit a classifier  $T^{(m)}(\mathbf{x})$  to the training data using weights  $w_i$ .

(b) Compute

$$err^{(m)} = \sum_{i=1}^n w_i \mathbb{I}(c_i \neq T^{(m)}(\mathbf{x}_i)) / \sum_{i=1}^n w_i.$$

(c) Compute

$$\alpha^{(m)} = \log \frac{1 - err^{(m)}}{err^{(m)}} + \log(K - 1).$$

(d) Set

$$w_i \leftarrow w_i \cdot \exp\left(\alpha^{(m)} \cdot \mathbb{I}(c_i \neq T^{(m)}(\mathbf{x}_i))\right), \quad i = 1, \dots, n.$$

(e) Re-normalize  $w_i$ .

3. Output

$$C(\mathbf{x}) = \arg \max_k \sum_{m=1}^M \alpha^{(m)} \cdot \mathbb{I}(T^{(m)}(\mathbf{x}) = k).$$

Questions?



# Resources

- Friedman, Jerome H. “Stochastic Gradient Boosting.” Statweb.stanford.edu, 26 Mar. 1999, [statweb.stanford.edu/~jhf/ftp/stobst.pdf](http://statweb.stanford.edu/~jhf/ftp/stobst.pdf).
- Shalev-Shwartz, Shai, and Shai Ben-David. Understanding Machine Learning: from Theory to Algorithms. Cambridge University Press, 2019.
- Zhu, Ji, et al. “Multi-Class AdaBoost.” Web.stanford.edu, 12 Jan. 2006, [web.stanford.edu/~hastie/Papers/samme.pdf](http://web.stanford.edu/~hastie/Papers/samme.pdf).



# Data

## Experiment:

- Keras.datasets, FashionMNIST, [https://keras.io/api/datasets/fashion\\_mnist/](https://keras.io/api/datasets/fashion_mnist/)
- Kaggle, Chinese MNIST, <https://www.kaggle.com/gpreda/chinese-mnist>
- Kaggle, Sign Language MNIST, <https://www.kaggle.com/datamunge/sign-language-mnist>

## Homework:

- UCI, Speaker Accent Recognition, <https://archive.ics.uci.edu/ml/datasets/Speaker+Accent+Recognition#>
- UCI, Deepfakes: Medical Image Tamper Detection, <https://archive.ics.uci.edu/ml/datasets/Deepfakes%3A+Medical+Image+Tamper+Detection>