
TOOLING ON MATLAB FOR ONLINE CONVEX OPTIMIZATION

Philip Naveen
University of Virginia
Charlottesville
philipnaveen@email.virginia.edu

ABSTRACT

This manuscript is merely a formal documentation of the purpose and details surrounding the online convex optimization toolbox (OCOBox) for MATLAB. The purpose of this toolbox is to provide a collection of algorithms that work under stochastic situations where traditional algorithmic theory does not fare so well. The toolbox encompasses a wide range of methods including Bayesian persuasion, bandit optimization, Blackwell approachability, boosting, game theory, projection-free algorithms, and regularization. In the future, we plan to extend OCOBox to interactive machine learning algorithms and develop a more robust GUI.

Keywords MATLAB · Online Convex Optimization · Stochastic Processes

1 Introduction

Optimization is a fundamental area of mathematics and computer science, with applications spanning a wide range of fields including machine learning, operations research, economics, engineering, and more. Optimization algorithms seek to find the best solution from a set of possible solutions, often by maximizing or minimizing a certain objective function. The diversity and complexity of optimization problems have led to the development of numerous algorithms, each with its unique strengths, weaknesses, and areas of applicability [7].

MATLAB, a high-level language and interactive environment for numerical computation, visualization, and programming, is widely used in both academic and industrial settings for developing and implementing optimization algorithms. However, researchers and practitioners often face challenges when it comes to testing and comparing different optimization methods due to the lack of a unified framework [5].

To address this challenge, we introduce the OCOBox, a comprehensive MATLAB toolbox that consolidates a wide array of optimization algorithms into a single, user-friendly interface. The toolbox is designed to facilitate the implementation, comparison, and analysis of various optimization techniques, making it an invaluable resource for researchers, educators, and practitioners alike.

The primary motivation behind the development of the OCOBox is to provide a versatile and accessible platform for exploring optimization algorithms. Specifically, the toolbox aims to achieve the following objectives:

1. **Comprehensive Coverage:** Include a diverse set of optimization algorithms covering various categories such as Bayesian persuasion, bandit optimization, gradient descent methods, boosting, game theory, and more. This breadth ensures that users can find suitable algorithms for a wide range of optimization problems.
2. **Ease of Use:** Offer a user-friendly collection of algorithms to save time from searching the internet and textbook for implementations.
3. **Comparative Analysis:** Facilitate the comparative analysis of different optimization algorithms by providing a standardized framework for evaluation. Users can easily switch between algorithms and compare their performance on various benchmark problems.
4. **Educational Value:** Serve as a valuable educational resource for teaching optimization concepts. The toolbox includes well-documented code and examples, making it an excellent tool for students and educators in optimization courses.

1.1 Overview of the Toolbox

The OCOBox is organized into several subdirectories, each corresponding to a specific category of optimization algorithms. This modular structure allows users to navigate the toolbox easily and focus on the algorithms relevant to their needs. The main categories and their respective algorithms are as follows:

- **Algorithmic Bayesian Persuasion:** Models and solves problems where a sender influences the beliefs and actions of a receiver through strategic information disclosure.
- **Bandit Optimization:** Addresses decision-making problems where an agent must balance exploration and exploitation.
- **Blackwell Approachability:** Applies game-theoretic concepts to ensure outcomes converge to a desired set.
- **Boosting:** Implements ensemble methods to improve the performance of weak learners.
- **Games and Duality:** Includes algorithms related to game theory and linear programming duality.
- **Gradient Descent:** Provides various methods for gradient-based optimization.
- **Learning Theory:** Covers foundational algorithms in statistical learning theory.
- **Momentum-Based Optimization:** Includes algorithms that accelerate gradient descent by considering past gradients.
- **Online Optimization:** Updates models incrementally as new data arrives.
- **Projection-Free Algorithms:** Utilizes methods that avoid computationally expensive projections.
- **Regularization:** Prevents overfitting by adding constraints or penalties.

Each algorithm in the toolbox is implemented as a MATLAB script or function, adhering to best practices in software development to ensure readability, modularity, and reusability. The toolbox includes detailed documentation and examples for each algorithm, guiding users through the process of setting up and running optimization tasks.

The OCOBox is applicable to a wide range of domains. In machine learning, it can be used to fine-tune models and improve performance through stochastic optimization and regularization techniques [8]. In operations research, it provides tools for solving complex decision-making problems. In economics, it offers methods for strategic information disclosure and game theory analysis. Furthermore, its educational value makes it an excellent resource for teaching optimization concepts in academic settings. The algorithms are almost exclusively based on textbooks on the OCO, and research conducted on algorithmic Bayesian persuasion [4, 8, 3].

2 Algorithmic Bayesian Persuasion

Bayesian persuasion involves a sender who aims to influence the receiver's actions by sharing information about an underlying state of the world. The sender's goal is to design an optimal signaling strategy that maximizes their payoff, considering the receiver's response to the information provided. This concept has applications in economics, political science, and beyond.

The Optimization Toolbox UI includes several MATLAB classes that implement different Bayesian persuasion strategies. These classes are:

- **BayesianPersuasion:** A general framework for Bayesian persuasion.
- **BayesianPersuasionArbitrary:** A strategy where the sender picks arbitrary hypotheses.
- **BayesianPersuasionMajority:** A strategy where the sender computes the majority hypothesis.
- **BayesianPersuasionRandomized:** A strategy where the sender picks random hypotheses.

3 BayesianPersuasion Class

Bayesian persuasion involves a sender who aims to influence the receiver's actions by sharing information about an underlying state of the world. The sender's goal is to design an optimal signaling strategy that maximizes their payoff, considering the receiver's response to the information provided. This concept has applications in economics, political science, and beyond.

The Optimization Toolbox UI includes several MATLAB classes that implement different Bayesian persuasion strategies. These classes are:

- **BayesianPersuasion**: A general framework for Bayesian persuasion.
- **BayesianPersuasionArbitrary**: A strategy where the sender picks arbitrary hypotheses.
- **BayesianPersuasionMajority**: A strategy where the sender computes the majority hypothesis.
- **BayesianPersuasionRandomized**: A strategy where the sender picks random hypotheses.

4 BayesianPersuasion Class

4.1 Overview

The `BayesianPersuasion` class provides a general framework for Bayesian persuasion. It defines the properties of the states, actions, sender's payoffs, receiver's payoffs, and the prior distribution over states. It also includes methods for finding the optimal signaling strategy and calculating the expected payoff for the sender.

4.2 Inputs and Outputs

- **Inputs:**
 - `states`: The possible states of the world.
 - `actions`: The actions available to the receiver.
 - `senderPayoff`: The sender's payoff matrix.
 - `receiverPayoff`: The receiver's payoff matrix.
 - `prior`: The prior distribution over the states.
- **Outputs:**
 - `sigma`: The optimal signaling strategy.
 - `rho`: The optimal response strategy.
 - `expectedPayoff`: The expected payoff for the sender given the strategies.

4.2.1 Algorithm

Algorithm 1 BayesianPersuasion: Finding Optimal Strategy

```

1: Initialize  $\sigma$  and  $\rho$  with random values.
2: for each state  $i$  do
3:   Normalize  $\sigma(i, :)$  to sum to 1.
4:   Normalize  $\rho(i, :)$  to sum to 1.
5: end for
6: return  $\sigma, \rho$ 

```

Algorithm 2 BayesianPersuasion: Calculating Expected Payoff

```

1: Initialize expectedPayoff to 0.
2: for each state  $i$  do
3:   for each action  $j$  do
4:     Update expectedPayoff using prior(i), sigma(i, j), senderPayoff(i, j), and rho(j, i).
5:   end for
6: end for
7: return expectedPayoff

```

5 BayesianPersuasionArbitrary Class

5.1 Overview

The `BayesianPersuasionArbitrary` class implements a strategy where the sender selects arbitrary hypotheses until no counter-examples are found. This approach iteratively refines the hypothesis space based on counter-examples.

5.2 Inputs and Outputs

- **Inputs:**
 - H : The concept class matrix.
 - P : The probability distribution over the feature space.
- **Outputs:**
 - h : The selected hypothesis.
 - x : The counter-example provided by the teacher.

5.2.1 Algorithm

Algorithm 3 BayesianPersuasionArbitrary: Arbitrary Learning

```

1: Initialize  $i$  to 1.
2: while true do
3:   Select an arbitrary hypothesis  $h$  from  $H$ .
4:   Obtain a counter-example  $x$  from the teacher.
5:   if  $x$  is empty then
6:     break
7:   end if
8:   Remove inconsistent hypotheses from  $H$ .
9:   Increment  $i$ .
10: end while
11: return  $h, x$ 

```

6 BayesianPersuasionMajority Class

6.1 Overview

The `BayesianPersuasionMajority` class implements a strategy where the sender computes the majority hypothesis. This method iteratively refines the hypothesis space based on majority votes and counter-examples.

6.2 Inputs and Outputs

- **Inputs:**
 - H : The concept class matrix.
 - P : The probability distribution over the feature space.
- **Outputs:**
 - h : The majority hypothesis.
 - x : The counter-example provided by the teacher.

6.2.1 Algorithm

Algorithm 4 BayesianPersuasionMajority: Majority Learning

```

1: while true do
2:   Compute the majority hypothesis  $h$ .
3:   Obtain a counter-example  $x$  from the teacher.
4:   if  $x$  is empty then
5:     break
6:   end if
7:   Remove inconsistent hypotheses from  $H$ .
8: end while
9: return  $h, x$ 

```

7 BayesianPersuasionRandomized Class

7.1 Overview

The `BayesianPersuasionRandomized` class implements a strategy where the sender randomly picks hypotheses. This approach iteratively refines the hypothesis space based on random selections and counter-examples.

7.2 Inputs and Outputs

- **Inputs:**
 - H : The concept class matrix.
 - P : The probability distribution over the feature space.
- **Outputs:**
 - h : The randomly selected hypothesis.
 - x : The counter-example provided by the teacher.

7.2.1 Algorithm

Algorithm 5 BayesianPersuasionRandomized: Randomized Learning

```

1: while true do
2:   Randomly select a hypothesis  $h$  from  $H$ .
3:   Obtain a counter-example  $x$  from the teacher.
4:   if  $x$  is empty then
5:     break
6:   end if
7:   Remove inconsistent hypotheses from  $H$ .
8: end while
9: return  $h, x$ 

```

8 Bandit Optimization

Bandit optimization algorithms are widely used in scenarios where decision-makers need to balance exploration and exploitation to maximize rewards. The Optimization Toolbox UI provides several MATLAB functions implementing bandit optimization algorithms for different scenarios.

9 banditLinearOptimization Function

9.1 Overview

The `banditLinearOptimization` function performs bandit optimization for linear functions. It optimizes the function f defined by A and b starting from x_0 using stepsize α for $maxIter$ iterations with perturbation parameter δ .

9.1.1 Algorithm

9.1.2 Inputs

- f - Function handle of the objective function.
- A - Matrix defining the linear function.
- b - Vector defining the linear function.
- x_0 - Initial guess.
- α - Step size.
- $maxIter$ - Maximum number of iterations.
- δ - Perturbation parameter.

Algorithm 6 banditLinearOptimization

```

1: procedure BANDITLINEAROPTIMIZATION( $f, A, b, x_0, \alpha, \maxIter, \delta$ )
2:   Inputs:  $f, A, b, x_0, \alpha, \maxIter, \delta$ 
3:   Outputs:  $x, history$ 
4:    $x \leftarrow x_0$ 
5:    $history \leftarrow \text{zeros}(\maxIter, 1)$ 
6:   for  $k = 1$  to  $\maxIter$  do
7:      $u \leftarrow \text{randn}(n, 1)$ 
8:      $u \leftarrow u/\text{norm}(u)$  ▷ Unit vector
9:      $fk\_plus \leftarrow f(x + \delta \cdot u, A, b)$ 
10:     $fk\_minus \leftarrow f(x - \delta \cdot u, A, b)$ 
11:     $gradEstimate \leftarrow (fk\_plus - fk\_minus)/(2 \cdot \delta) \cdot u$ 
12:     $x \leftarrow x - \alpha \cdot gradEstimate$ 
13:     $history(k) \leftarrow f(x, A, b)$ 
14:   end for
15: end procedure

```

9.1.3 Outputs

- x - Optimized parameters.
- $history$ - History of function values.

10 banditOptimization Function**10.1 Overview**

The banditOptimization function performs bandit optimization using gradient estimates. It optimizes the function f starting from x_0 using stepsize α for \maxIter iterations with perturbation parameter δ .

10.1.1 Algorithm**Algorithm 7** banditOptimization

```

1: procedure BANDITOPTIMIZATION( $f, x_0, \alpha, \maxIter, \delta$ )
2:   Inputs:  $f, x_0, \alpha, \maxIter, \delta$ 
3:   Outputs:  $x, history$ 
4:    $x \leftarrow x_0$ 
5:    $history \leftarrow \text{zeros}(\maxIter, 1)$ 
6:   for  $k = 1$  to  $\maxIter$  do
7:      $u \leftarrow \text{randn}(n, 1)$ 
8:      $u \leftarrow u/\text{norm}(u)$  ▷ Unit vector
9:      $fk\_plus \leftarrow f(x + \delta \cdot u)$ 
10:     $fk\_minus \leftarrow f(x - \delta \cdot u)$ 
11:     $gradEstimate \leftarrow (fk\_plus - fk\_minus)/(2 \cdot \delta) \cdot u$ 
12:     $x \leftarrow x - \alpha \cdot gradEstimate$ 
13:     $history(k) \leftarrow f(x)$ 
14:   end for
16: end procedure

```

10.1.2 Inputs

- f - Function handle of the objective function.
- x_0 - Initial guess.
- α - Step size.
- \maxIter - Maximum number of iterations.
- δ - Perturbation parameter.

10.1.3 Outputs

- x - Optimized parameters.
- $history$ - History of function values.

11 exp3 Function

11.1 Overview

The `exp3` function performs the EXP3 algorithm for the multi-armed bandit problem. It runs the EXP3 algorithm for T rounds with K arms and exploration parameter γ .

11.1.1 Algorithm

Algorithm 8 `exp3`

```

1: procedure EXP3( $T, K, \gamma$ )
2:   Inputs:  $T, K, \gamma$ 
3:   Outputs:  $weights, history$ 
4:    $weights \leftarrow \text{ones}(K, 1)$ 
5:    $history \leftarrow \text{zeros}(T, 2)$  ▷ Column 1: chosen arm, Column 2: received reward
6:   for  $t = 1$  to  $T$  do
7:      $probs \leftarrow (1 - \gamma) \cdot (weights / \sum(weights)) + (\gamma / K)$ 
8:      $arm \leftarrow \text{randSample}(1 : K, 1, \text{true}, probs)$ 
9:      $reward \leftarrow \text{banditReward}(arm, t)$ 
10:     $estimatedReward \leftarrow reward / probs(arm)$ 
11:     $weights(arm) \leftarrow weights(arm) \cdot \exp(\gamma \cdot estimatedReward / K)$ 
12:     $history(t, :) \leftarrow [arm, reward]$ 
13:   end for
14: end procedure

```

11.1.2 Inputs

- T - Number of rounds.
- K - Number of arms.
- γ - Exploration parameter.

11.1.3 Outputs

- $weights$ - Final weights of the arms.
- $history$ - History of chosen arms and rewards.

12 exp3WithEstimators Function

12.1 Overview

The `exp3WithEstimators` function performs the EXP3 algorithm using unbiased estimators. It runs the EXP3 algorithm for T rounds with K arms and exploration parameter γ using unbiased estimators for the rewards.

12.1.1 Algorithm

This function shares the same algorithmic structure as the `exp3` function, with the difference lying in the reward calculation, which employs unbiased estimators.

12.1.2 Inputs

- T - Number of rounds.
- K - Number of arms.

Algorithm 9 exp3WithEstimators

```

1: procedure EXP3WITHESTIMATORS( $T, K, \gamma$ )
2:   Inputs:  $T, K, \gamma$ 
3:   Outputs:  $weights, history$ 
4:    $weights \leftarrow \text{ones}(K, 1)$ 
5:    $history \leftarrow \text{zeros}(T, 2)$  ▷ Column 1: chosen arm, Column 2: received reward
6:   for  $t = 1$  to  $T$  do
7:      $probs \leftarrow (1 - \gamma) \cdot (weights / \sum(weights)) + (\gamma / K)$ 
8:      $arm \leftarrow \text{randSample}(1 : K, 1, \text{true}, probs)$ 
9:      $reward \leftarrow \text{banditRewardWithEstimators}(arm, t)$ 
10:     $estimatedReward \leftarrow reward / probs(arm)$ 
11:     $weights(arm) \leftarrow weights(arm) \cdot \exp(\gamma \cdot estimatedReward / K)$ 
12:     $history(t, :) \leftarrow [arm, reward]$ 
13:   end for
14: end procedure

```

- γ - Exploration parameter.

12.1.3 Outputs

- $weights$ - Final weights of the arms.
- $history$ - History of chosen arms and rewards.

13 Blackwell Approachability

The approachability function computes the approachability solution for a given matrix A and vector b .

13.0.1 Algorithm**Algorithm 10** approachability

```

1: procedure APPROACHABILITY( $A, b$ )
2:   Inputs:  $A, b$ 
3:   Output:  $x_{\text{opt}}$ 
4:    $n \leftarrow \text{size}(A, 2)$ 
5:   cvx_begin quiet
6:   variable  $x^{(n)}$ 
7:   minimize  $\|A \cdot x - b\|$ 
8:   cvx_end
9:    $x_{\text{opt}} \leftarrow x$ 
10: end procedure

```

13.0.2 Inputs

- A - Matrix A .
- b - Vector b .

13.0.3 Outputs

- x_{opt} - Approachability solution.

13.1 polarCones Function**13.2 Overview**

The polarCones function computes the polar cones for a given matrix A and vector b .

13.2.1 Algorithm

Algorithm 11 polarCones

```

1: procedure POLARCONES( $A, b$ )
2:   Inputs:  $A, b$ 
3:   Outputs:  $cones, dualCones$ 
4:    $[m, n] \leftarrow \text{size}(A)$ 
5:    $cones \leftarrow \text{cell}(n, 1)$ 
6:    $dualCones \leftarrow \text{cell}(m, 1)$ 
7:   for  $i = 1$  to  $n$  do
8:     cvx_begin quiet
9:     variable  $x(n)$ 
10:    minimize  $\|A \cdot x - b\|$ 
11:    subject to  $\|x(i+1 : \text{end})\|_2 \leq x(i)$ 
12:    cvx_end
13:     $cones\{i\} \leftarrow x$ 
14:  end for
15:  for  $j = 1$  to  $m$  do
16:    cvx_begin quiet
17:    variable  $y(m)$ 
18:    maximize  $y' \cdot b$ 
19:    subject to  $y' \cdot A \leq 0$ 
20:     $y(j) = 1$ 
21:    cvx_end
22:     $dualCones\{j\} \leftarrow y$ 
23:  end for
24: end procedure

```

Inputs

- A - Matrix A .
- b - Vector b .

13.2.2 Outputs

- $cones$ - Polar cones.
- $dualCones$ - Dual polar cones.

14 Boosting

The adaboost function implements the AdaBoost algorithm, which trains an ensemble of base learners using the given base learning algorithm.

14.0.1 Algorithm

14.0.2 Inputs

- `baseLearner` - A function handle for the base learning algorithm with signature $@(X, y, weights)$.
- `X_train` - Training data features.
- `y_train` - Training data labels.
- `numRounds` - Number of boosting rounds.

14.0.3 Outputs

- `models` - Cell array of trained base learners.
- `alphas` - Weights for each base learner.

Algorithm 12 AdaBoost

```

1: procedure ADABOOST(baseLearner, X_train, y_train, numRounds)
2:   n ← size(X_train, 1)
3:   weights ← ones(n, 1)/n                                     ▷ Initialize weights
4:   models ← cell(numRounds, 1)
5:   alphas ← zeros(numRounds, 1)
6:   for t ← 1 to numRounds do
7:     models{t} ← baseLearner(X_train, y_train, weights)         ▷ Train base learner
8:     y_pred ← predict(models{t}, X_train)                       ▷ Predict with base learner
9:      $\epsilon_t \leftarrow \sum_{i=1}^n (\text{weights}(i) \cdot (y\_pred(i) \neq y\_train(i)))$    ▷ Compute weighted error
10:     $\alpha_t \leftarrow 0.5 \cdot \log\left(\frac{1-\epsilon_t}{\epsilon_t}\right)$            ▷ Compute alpha
11:    weights ← updateWeights(weights,  $\alpha_t$ , y_train, y_pred)   ▷ Update weights
12:    models{t} ← models{t}                                       ▷ Store model
13:    alphas(t) ←  $\alpha_t$                                            ▷ Store alpha
14:   end for
15:   return models, alphas
16: end procedure

```

14.0.4 boostingByOCO Function

The `boostingByOCO` function implements boosting using Online Convex Optimization (OCO).

14.0.5 Algorithm**Algorithm 13** Boosting by OCO

```

1: procedure BOOSTINGBYOCO(baseLearner, X_train, y_train, numRounds)
2:   n ← size(X_train, 1)
3:   weights ← ones(n, 1)/n                                     ▷ Initialize weights
4:   models ← cell(numRounds, 1)
5:   for t ← 1 to numRounds do
6:     models{t} ← baseLearner(X_train, y_train, weights)         ▷ Train base learner
7:     y_pred ← predict(models{t}, X_train)                       ▷ Predict with base learner
8:      $\epsilon_t \leftarrow \sum_{i=1}^n (\text{weights}(i) \cdot (y\_pred(i) \neq y\_train(i)))$    ▷ Compute weighted error
9:      $\alpha_t \leftarrow 0.5 \cdot \log\left(\frac{1-\epsilon_t}{\epsilon_t}\right)$            ▷ Compute alpha
10:    weights ← updateWeights(weights,  $\alpha_t$ , y_train, y_pred)   ▷ Update weights
11:    models{t} ← models{t}                                       ▷ Store model
12:   end for
13:   return weights, models
14: end procedure

```

14.0.6 Inputs

- `baseLearner` - A function handle for the base learning algorithm with signature `@(X, y, weights)`.
- `X_train` - Training data features.
- `y_train` - Training data labels.
- `numRounds` - Number of boosting rounds.

14.0.7 Outputs

- `models` - Cell array of trained base learners.
- `alphas` - Weights for each base learner.

15 AdaBoost Algorithms

15.1 adaboost Function

Description The `adaboost` function implements the AdaBoost algorithm, which trains an ensemble of base learners using the given base learning algorithm.

Algorithm

Algorithm 14 adaboost

```

1: Initialize weights:  $w_i = \frac{1}{n}$  for all  $i = 1, \dots, n$ 
2: for  $t = 1$  to numRounds do
3:   Train base learner with current weights:  $h_t = \text{baseLearner}(X, y, w)$ 
4:   Predict with base learner:  $\hat{y} = h_t(X)$ 
5:   Compute weighted error:  $\epsilon_t = \sum_{i=1}^n w_i \mathbf{1}(\hat{y}_i \neq y_i)$ 
6:   Compute alpha:  $\alpha_t = 0.5 \log\left(\frac{1-\epsilon_t}{\epsilon_t}\right)$ 
7:   Update weights:  $w_i \leftarrow w_i \exp(-\alpha_t y_i \hat{y}_i)$ 
8:   Normalize weights:  $w \leftarrow \frac{w}{\sum w}$ 
9: end for
10: Return models and alphas:  $\{h_t\}_{t=1}^{\text{numRounds}}, \{\alpha_t\}_{t=1}^{\text{numRounds}}$ 

```

Inputs

- `baseLearner` - A function handle for the base learning algorithm with signature `@(X, y, weights)`.
- `X_train` - Training data features.
- `y_train` - Training data labels.
- `numRounds` - Number of boosting rounds.

Outputs

- `models` - Cell array of trained base learners.
- `alphas` - Weights for each base learner.

15.2 boostingByOCO Function

Description The `boostingByOCO` function implements boosting using Online Convex Optimization (OCO).

Algorithm

Algorithm 15 boostingByOCO

```

1: Initialize weights:  $w_i = \frac{1}{n}$  for all  $i = 1, \dots, n$ 
2: for  $t = 1$  to numRounds do
3:   Train base learner with current weights:  $h_t = \text{baseLearner}(X, y, w)$ 
4:   Predict with base learner:  $\hat{y} = h_t(X)$ 
5:   Compute weighted error:  $\epsilon_t = \sum_{i=1}^n w_i \mathbf{1}(\hat{y}_i \neq y_i)$ 
6:   Compute alpha:  $\alpha_t = 0.5 \log\left(\frac{1-\epsilon_t}{\epsilon_t}\right)$ 
7:   Update weights:  $w_i \leftarrow w_i \exp(-\alpha_t y_i \hat{y}_i)$ 
8:   Normalize weights:  $w \leftarrow \frac{w}{\sum w}$ 
9: end for
10: Return weights and models:  $w, \{h_t\}_{t=1}^{\text{numRounds}}$ 

```

Inputs

- `baseLearner` - A function handle for the base learning algorithm with signature `@(X, y, weights)`.
- `X_train` - Training data features.
- `y_train` - Training data labels.
- `numRounds` - Number of boosting rounds.

Outputs

- `weights` - Final weights for each base learner.
- `models` - Cell array of trained base learners.

15.3 contextualLearningModel Function

Description The `contextualLearningModel` function implements a contextual learning model using boosting.

Algorithm

Algorithm 16 contextualLearningModel

- 1: Initialize weights: $w_i = \frac{1}{n}$ for all $i = 1, \dots, n$
 - 2: **for** $t = 1$ to `numRounds` **do**
 - 3: Train base learner with current weights and context: $h_t = \text{baseLearner}(X, y, w, \text{contextFeature})$
 - 4: Predict with base learner: $\hat{y} = h_t(X)$
 - 5: Compute weighted error: $\epsilon_t = \sum_{i=1}^n w_i \mathbf{1}(\hat{y}_i \neq y_i)$
 - 6: Compute alpha: $\alpha_t = 0.5 \log\left(\frac{1-\epsilon_t}{\epsilon_t}\right)$
 - 7: Update weights: $w_i \leftarrow w_i \exp(-\alpha_t y_i \hat{y}_i)$
 - 8: Normalize weights: $w \leftarrow \frac{w}{\sum w}$
 - 9: **end for**
 - 10: Return model: $\{\text{contextFeature}, \{h_t\}_{t=1}^{\text{numRounds}}, \{\alpha_t\}_{t=1}^{\text{numRounds}}\}$
-

Inputs

- `baseLearner` - A function handle for the base learning algorithm with signature `@(X, y, context)`.
- `contextFeature` - Contextual feature vector.
- `X_train` - Training data features.
- `y_train` - Training data labels.
- `numRounds` - Number of boosting rounds.

Outputs

- `model` - Trained contextual learning model.

15.4 onlineBoosting Function

Description The `onlineBoosting` function implements the Online Boosting algorithm.

Algorithm

Inputs

- `baseLearner` - A function handle for the base learning algorithm with signature `@(X, y)`.
- `X_train` - Training data features.
- `y_train` - Training data labels.
- `numRounds` - Number of boosting rounds.

Algorithm 17 onlineBoosting

```

1: Initialize weights:  $w_i = \frac{1}{n}$  for all  $i = 1, \dots, n$ 
2: for  $t = 1$  to numRounds do
3:   Train base learner with current weights:  $h_t = \text{baseLearner}(X, y)$ 
4:   Predict with base learner:  $\hat{y} = h_t(X)$ 
5:   Compute weighted error:  $\epsilon_t = \sum_{i=1}^n w_i \mathbf{1}(\hat{y}_i \neq y_i)$ 
6:   Compute alpha:  $\alpha_t = 0.5 \log\left(\frac{1-\epsilon_t}{\epsilon_t}\right)$ 
7:   Update weights:  $w_i \leftarrow w_i \exp(-\alpha_t y_i \hat{y}_i)$ 
8:   Normalize weights:  $w \leftarrow \frac{w}{\sum w}$ 
9: end for
10: Return models and alphas:  $\{h_t\}_{t=1}^{\text{numRounds}}, \{\alpha_t\}_{t=1}^{\text{numRounds}}$ 

```

Outputs

- models - Cell array of trained base learners.
- alphas - Weights for each base learner.

15.5 predictBoostedModel Function

Description The predictBoostedModel function predicts using a boosted model.

Algorithm

Algorithm 18 predictBoostedModel

```

1: Initialize aggregated predictions:  $F(x) = 0$  for all  $x$ 
2: for  $t = 1$  to numRounds do
3:   Predict with base learner:  $\hat{y}_t = h_t(X)$ 
4:   Aggregate predictions:  $F(x) \leftarrow F(x) + \alpha_t \hat{y}_t$ 
5: end for
6: Return final prediction:  $\hat{y} = \text{sign}(F(x))$ 

```

Inputs

- model - A struct containing the trained boosted model with fields 'models' and 'alphas'.
- X - Input data features.

Outputs

- y_pred - Predicted labels.

16 Approximating Linear Programs**16.1 approximatingLinearPrograms Function**

Description The approximatingLinearPrograms function solves a linear programming problem approximately by optimizing the function $c'x$ subject to $A * x \leq b$ using an iterative method starting from an initial guess x_0 .

Algorithm

Inputs

- A - Constraint matrix.
- b - Constraint vector.

Algorithm 19 approximatingLinearPrograms

```

1: Initialize  $x = x_0$ 
2: Initialize history array to store objective values
3: for  $k = 1$  to  $\text{maxIter}$  do
4:   Calculate subgradient:  $g = A(\text{idx}, :)'$  where  $\text{idx} = \arg \max(A * x - b)$ 
5:   Compute stepsize:  $\alpha = \frac{1}{\sqrt{k}}$ 
6:   Update  $x$ :  $x = x - \alpha * (g + c)$ 
7:   Project  $x$  onto feasible set:  $x = \text{projectionOntoFeasibleSet}(x, A, b)$ 
8:   Calculate objective value:  $\text{history}(k) = c' * x$   $\text{norm}(g) < \text{tol}$ 
9:   Break and store history up to iteration  $k$ 
10: end for
11: Return solution  $x$  and history of objective values

```

- c - Cost vector.
- x_0 - Initial guess.
- maxIter - Maximum number of iterations.
- tol - Tolerance for convergence.

Outputs

- x - Solution to the linear programming problem.
- history - History of objective function values.

16.2 projectionOntoFeasibleSet Function

Description The `projectionOntoFeasibleSet` function projects the vector x onto the feasible set defined by $A * x \leq b$.

Algorithm

Algorithm 20 projectionOntoFeasibleSet

```

1: Solve quadratic program:  $x = \text{quadprog}(\text{eye}(\text{length}(x)), -x, A, b, \text{options})$ 
2: Return projected  $x$ 

```

Inputs

- x - Vector to be projected.
- A - Constraint matrix.
- b - Constraint vector.

Outputs

- x - Projected vector.

17 Linear Programming Duality

17.1 linearProgrammingDuality Function

Description The `linearProgrammingDuality` function solves the primal and dual linear programming problems.

Algorithm

Inputs

Algorithm 21 linearProgrammingDuality

- 1: Solve primal problem: `[primal, , , , lambda] = linprog(c, A, b, [], [], [], [], options)`
 - 2: Extract dual solution: `dual = lambda.ineqlin`
 - 3: Return primal and dual solutions
-

- `c` - Cost vector.
- `A` - Constraint matrix.
- `b` - Constraint vector.

Outputs

- `primal` - Solution to the primal problem.
- `dual` - Solution to the dual problem.

18 Von Neumann Minimax Theorem

18.1 vonNeumannTheorem Function

Description The `vonNeumannTheorem` function verifies the von Neumann Minimax Theorem by returning the value of the game using von Neumann's theorem.

Algorithm

Algorithm 22 vonNeumannTheorem

- 1: Solve zero-sum game: `[, , value] = zeroSumGame(A)`
 - 2: Return value of the game
-

Inputs

- `A` - Payoff matrix.

Outputs

- `value` - Value of the game.

19 Zero-Sum Game

19.1 zeroSumGame Function

Description The `zeroSumGame` function solves a zero-sum game using linear programming.

Algorithm

Inputs

- `A` - Payoff matrix.

Outputs

- `x` - Mixed strategy for player 1.
- `y` - Mixed strategy for player 2.
- `value` - Value of the game.

Algorithm 23 zeroSumGame

-
- 1: Define sizes: $[m, n] = \text{size}(A)$
 - 2: Solve for player 1:
 - Define problem: $f = [-\mathbf{1}_m; 1]$
 - Define constraints: $A1 = [A', -\mathbf{1}_n], b1 = \mathbf{0}_n$
 - Define equality: $Aeq1 = [\mathbf{1}'_m, 0], beq1 = 1$
 - Define bounds: $lb1 = [\mathbf{0}_m; -\infty]$
 - Solve: $x1 = \text{linprog}(f, A1, b1, Aeq1, beq1, lb1)$
 - 3: Solve for player 2:
 - Define problem: $g = [\mathbf{1}_n; -1]$
 - Define constraints: $A2 = [-A, \mathbf{1}_m], b2 = \mathbf{0}_m$
 - Define equality: $Aeq2 = [\mathbf{1}'_n, 0], beq2 = 1$
 - Define bounds: $lb2 = [\mathbf{0}_n; -\infty]$
 - Solve: $y2 = \text{linprog}(g, A2, b2, Aeq2, beq2, lb2)$
 - 4: Extract strategies and value:
 - $x = x1(1 : m)$
 - $y = y2(1 : n)$
 - $value = x1(end)$
 - 5: Return strategies x, y and game value
-

20 Constrained Gradient Descent

20.1 constrainedGradientDescent Function

Description The `constrainedGradientDescent` function performs constrained gradient descent optimization. It optimizes the function f with its gradient `gradf` starting from $x0$ using stepsize α for `maxIter` iterations, subject to the constraint specified by `constraintFunc`.

Algorithm

Algorithm 24 constrainedGradientDescent

-
- 1: Initialize $x = x0$
 - 2: Initialize history array to store function values
 - 3: **for** $k = 1$ to `maxIter` **do**
 - 4: Update x : $x = x - \alpha * \text{gradf}(x)$
 - 5: Project x onto the constraint set: $x = \text{constraintFunc}(x)$
 - 6: Store function value: `history(k) = f(x)`
 - 7: **end for**
 - 8: Return optimized parameters x and history of function values
-

Inputs

- `f` - Function handle of the objective function.
- `gradf` - Function handle of the gradient of f .
- `x0` - Initial guess.
- `alpha` - Step size.
- `maxIter` - Maximum number of iterations.
- `constraintFunc` - Function handle of the constraint projection function.

Outputs

- `x` - Optimized parameters.
- `history` - History of function values.

21 Gradient Descent

21.1 `gradientDescent` Function

Description The `gradientDescent` function performs gradient descent optimization. It optimizes the function f with its gradient `gradf` starting from x_0 using stepsize α for `maxIter` iterations.

Algorithm

Algorithm 25 `gradientDescent`

- 1: Initialize $x = x_0$
 - 2: Initialize history array to store function values
 - 3: **for** $k = 1$ to `maxIter` **do**
 - 4: Update x : $x = x - \alpha * \text{gradf}(x)$
 - 5: Store function value: `history(k) = f(x)`
 - 6: **end for**
 - 7: Return optimized parameters x and history of function values
-

Inputs

- `f` - Function handle of the objective function.
- `gradf` - Function handle of the gradient of f .
- `x0` - Initial guess.
- `alpha` - Step size.
- `maxIter` - Maximum number of iterations.

Outputs

- `x` - Optimized parameters.
- `history` - History of function values.

22 Distance to Optimality

22.1 `measureDistanceToOptimality` Function

Description The `measureDistanceToOptimality` function measures the Euclidean distance to the optimal point. It computes the Euclidean distance between the current point x and the optimal point x_{opt} .

Algorithm

Algorithm 26 `measureDistanceToOptimality`

- 1: Compute Euclidean distance: `distance = norm(x - x_opt)`
 - 2: Return distance
-

Inputs

- `x` - Current point.
- `x_opt` - Optimal point.

Outputs

- distance - Euclidean distance to the optimal point.

23 Polyak Stepsize

23.1 polyakStepsize Function

Description The `polyakStepsize` function computes the Polyak stepsize for gradient descent. It calculates the stepsize for the function f at point x given the optimal value f_{opt} .

Algorithm

Algorithm 27 `polyakStepsize`

- 1: Compute gradient: `gradf = gradient(f, x)`
 - 2: Compute Polyak stepsize: $\alpha = \frac{f(x) - f_{\text{opt}}}{\|\text{gradf}(x)\|^2}$
 - 3: Return stepsize α
-

Inputs

- `f` - Function handle of the objective function.
- `x` - Current point.
- `f_opt` - Optimal value of the function.

Outputs

- `alpha` - Polyak stepsize.

24 SVM Training

24.1 svmTraining Function

Description The `svmTraining` function trains a Support Vector Machine using gradient descent. It trains an SVM with regularization parameter λ using gradient descent for `maxIter` iterations and stepsize α .

Algorithm

Algorithm 28 `svmTraining`

- 1: Initialize weights w and bias b to zero
 - 2: Initialize history array to store loss values
 - 3: **for** `iter = 1 to maxIter` **do**
 - 4: **for** `i = 1 to m` **do** `y(i) * (X(i, :) * w + b) < 1`
 - 5: Update w : $w = w - \alpha * (\lambda * w - y(i) * X(i, :))'$
 - 6: Update b : $b = b + \alpha * y(i)$
 - 7: Update w : $w = w - \alpha * \lambda * w$
 - 8: **end for**
 - 9: Compute hinge loss: $\text{loss} = \text{sum}(\max(0, 1 - y .* (X * w + b))) / m + (\text{lambda} / 2) * \text{norm}(w)^2$
 - 10: Store loss: `history(iter) = loss`
 - 11: **end for**
 - 12: Return weights w , bias b , and history of loss values
-

Inputs

- X - Matrix of training examples (each row is a training example).
- y - Vector of labels (+1 or -1).
- λ - Regularization parameter.
- `maxIter` - Maximum number of iterations.
- `alpha` - Step size.

Outputs

- w - Weight vector.
- b - Bias term.
- `history` - History of loss values.

25 Generalization and Learnability

25.1 `generalizationAndLearnability` Function

Description The `generalizationAndLearnability` function evaluates the generalization error of a model with different complexity levels. It trains the model on the training data and computes the training and validation error for different complexity levels.

Algorithm

Algorithm 29 Generalization and Learnability Evaluation

```

1: procedure GENERALIZATIONANDLEARNABILITY(model,  $X_{\text{train}}$ ,  $y_{\text{train}}$ ,  $X_{\text{val}}$ ,  $y_{\text{val}}$ , complexity)
2:   trainError  $\leftarrow$  zeros(length(complexity), 1)
3:   valError  $\leftarrow$  zeros(length(complexity), 1)
4:   for  $i = 1$  to length(complexity) do
5:     Train the model with current complexity level:  $y_{\text{train\_pred}} \leftarrow \text{model}(X_{\text{train}}, y_{\text{train}}, X_{\text{train}}, \text{complexity}(i))$ 
6:     Compute training error:  $\text{trainError}(i) \leftarrow \text{mean}(y_{\text{train\_pred}} \neq y_{\text{train}})$ 
7:     Predict on validation data:  $y_{\text{val\_pred}} \leftarrow \text{model}(X_{\text{train}}, y_{\text{train}}, X_{\text{val}}, \text{complexity}(i))$ 
8:     Compute validation error:  $\text{valError}(i) \leftarrow \text{mean}(y_{\text{val\_pred}} \neq y_{\text{val}})$ 
9:   end for
10:  return trainError, valError
11: end procedure

```

Inputs

- `model` - A function handle for the learning model with signature `@(X, y, X_val, complexity)`.
- X_{train} - Training data features.
- y_{train} - Training data labels.
- X_{val} - Validation data features.
- y_{val} - Validation data labels.
- `complexity` - Vector of complexity levels.

Outputs

- `trainError` - Training error for each complexity level.
- `valError` - Validation error for each complexity level.

26 Statistical Learning Theory

26.1 statisticalLearningTheory Function

Description The `statisticalLearningTheory` function evaluates a model using training and testing data. It trains the model on the training data and computes the training and testing error.

Algorithm

Algorithm 30 Statistical Learning Theory Evaluation

```

1: procedure STATISTICALLEARNINGTHEORY(model,  $X_{\text{train}}$ ,  $y_{\text{train}}$ ,  $X_{\text{test}}$ ,  $y_{\text{test}}$ )
2:   Train the model:  $y_{\text{train\_pred}} \leftarrow \text{model}(X_{\text{train}}, y_{\text{train}}, X_{\text{train}})$ 
3:   Compute training error:  $\text{trainError} \leftarrow \text{mean}(y_{\text{train\_pred}} \neq y_{\text{train}})$ 
4:   Predict on test data:  $y_{\text{test\_pred}} \leftarrow \text{model}(X_{\text{train}}, y_{\text{train}}, X_{\text{test}})$ 
5:   Compute testing error:  $\text{testError} \leftarrow \text{mean}(y_{\text{test\_pred}} \neq y_{\text{test}})$ 
6:   return trainError, testError
7: end procedure

```

Inputs

- `model` - A function handle for the learning model with signature `@(X, y, X_val)`.
- X_{train} - Training data features.
- y_{train} - Training data labels.
- X_{test} - Testing data features.
- y_{test} - Testing data labels.

Outputs

- `trainError` - Training error.
- `testError` - Testing error.

27 AdaDelta Optimization

27.1 adadelta Function

Description The `adadelta` function implements the AdaDelta optimization algorithm. It performs optimization of the objective function f using the AdaDelta algorithm.

Algorithm

Algorithm 31 AdaDelta Optimization Algorithm

```

1: procedure ADADELTA( $f$ , grad_f,  $x_0$ , options)
2:   if nargin < 4 then
3:     options ← struct()
4:   end if
5:   rho ← 0.9 ▷ decay rate
6:   epsilon ← 1e - 8 ▷ small constant for numerical stability
7:   maxIter ← 1000 ▷ maximum number of iterations
8:   if isfield(options, 'rho') then
9:     rho ← options.rho
10:  end if
11:  if isfield(options, 'epsilon') then
12:    epsilon ← options.epsilon
13:  end if
14:  if isfield(options, 'maxIter') then
15:    maxIter ← options.maxIter
16:  end if
17:   $x \leftarrow x_0$  ▷ initialize
18:   $E_{\delta x} \leftarrow \text{zeros}(\text{size}(x))$  ▷ accumulated squared deltas
19:   $E_{\delta x_t} \leftarrow \text{zeros}(\text{size}(x))$  ▷ accumulated squared deltas (time varying)
20:  for  $t = 1$  to maxIter do ▷ optimization loop
21:    grad ← grad_f( $x$ ) ▷ compute gradient
22:    grad_sq ← grad.2 ▷ compute squared gradient
23:     $E_{\delta x} \leftarrow \text{rho} \times E_{\delta x} + (1 - \text{rho}) \times \text{grad\_sq}$  ▷ compute exponentially decaying average of squared deltas
24:     $\text{RMS}_{\delta x} \leftarrow \sqrt{E_{\delta x} + \text{epsilon}}$  ▷ compute RMS of squared deltas
25:     $\text{RMS}_{\delta x_t} \leftarrow \sqrt{E_{\delta x_t} + \text{epsilon}}$  ▷ compute RMS of squared deltas (time varying)
26:     $\Delta x \leftarrow - \left( \frac{\text{RMS}_{\delta x_t}}{\text{RMS}_{\delta x}} \right) \times \text{grad}$  ▷ update parameters
27:     $x \leftarrow x + \Delta x$  ▷ compute exponentially decaying average of squared deltas (time varying)
28:     $E_{\delta x_t} \leftarrow \text{rho} \times E_{\delta x_t} + (1 - \text{rho}) \times (\Delta x)^2$ 
29:  end for
30:   $f_{\text{opt}} \leftarrow f(x)$  ▷ compute optimal function value
31:  return  $x_{\text{opt}}, f_{\text{opt}}$ 
32: end procedure

```

Inputs

- f - Objective function.
- grad_f - Gradient of the objective function.
- x_0 - Initial point.
- options - Optimization options (optional).

Outputs

- x_{opt} - Optimal point.
- f_{opt} - Optimal function value.

28 Adagrad Optimization

28.1 adagrad Function

Description The adagrad function implements the Adagrad optimization algorithm. It performs optimization of the objective function f using the Adagrad algorithm.

Algorithm

Algorithm 32 Adagrad Optimization Algorithm

```

1: procedure ADAGRAD( $f$ ,  $\text{grad}_f$ ,  $x_0$ ,  $\text{options}$ )
2:   if  $\text{nargin} < 4$  then
3:      $\text{options} \leftarrow \text{struct}()$ 
4:   end if
5:    $\alpha \leftarrow 0.01$  ▷ learning rate
6:    $\epsilon \leftarrow 1e - 8$  ▷ small constant for numerical stability
7:    $\text{maxIter} \leftarrow 1000$  ▷ maximum number of iterations
8:   if  $\text{isfield}(\text{options}, 'alpha')$  then
9:      $\alpha \leftarrow \text{options.alpha}$ 
10:  end if
11:  if  $\text{isfield}(\text{options}, 'epsilon')$  then
12:     $\epsilon \leftarrow \text{options.epsilon}$ 
13:  end if
14:  if  $\text{isfield}(\text{options}, 'maxIter')$  then
15:     $\text{maxIter} \leftarrow \text{options.maxIter}$ 
16:  end if
17:   $x \leftarrow x_0$  ▷ initialize
18:   $G \leftarrow \text{zeros}(\text{size}(x))$  ▷ sum of squares of past gradients
19:  for  $t = 1$  to  $\text{maxIter}$  do ▷ optimization loop
20:     $\text{grad} \leftarrow \text{grad}_f(x)$  ▷ compute gradient
21:     $G \leftarrow G + \text{grad}^2$  ▷ accumulate squared gradient
22:     $x \leftarrow x - \alpha \times \text{grad} / (\sqrt{G} + \epsilon)$  ▷ update parameters
23:  end for
24:   $f_{\text{opt}} \leftarrow f(x)$  ▷ compute optimal function value
25:  return  $x_{\text{opt}}, f_{\text{opt}}$ 
26: end procedure

```

Inputs

- f - Objective function.
- grad_f - Gradient of the objective function.
- x_0 - Initial point.
- options - Optimization options (optional).

Outputs

- x_{opt} - Optimal point.
- f_{opt} - Optimal function value.

29 Adam Optimization

29.1 adam Function

Description The adam function implements the Adam optimization algorithm. It performs optimization of the objective function f using the Adam algorithm.

Algorithm

Algorithm 33 Adam Optimization Algorithm

```

1: procedure ADAM( $f, \text{grad}_f, x_0, \text{options}$ )
2:   if nargin < 4 then
3:     options  $\leftarrow$  struct()
4:   end if
5:   alpha  $\leftarrow$  0.001 ▷ learning rate
6:   beta1  $\leftarrow$  0.9 ▷ decay rate for 1st moment estimate
7:   beta2  $\leftarrow$  0.999 ▷ decay rate for 2nd moment estimate
8:   epsilon  $\leftarrow$   $1e - 8$  ▷ small constant for numerical stability
9:   maxIter  $\leftarrow$  1000 ▷ maximum number of iterations
10:  if isfield(options, 'alpha') then
11:    alpha  $\leftarrow$  options.alpha
12:  end if
13:  if isfield(options, 'beta1') then
14:    beta1  $\leftarrow$  options.beta1
15:  end if
16:  if isfield(options, 'beta2') then
17:    beta2  $\leftarrow$  options.beta2
18:  end if
19:  if isfield(options, 'epsilon') then
20:    epsilon  $\leftarrow$  options.epsilon
21:  end if
22:  if isfield(options, 'maxIter') then
23:    maxIter  $\leftarrow$  options.maxIter
24:  end if
25:   $x \leftarrow x_0$  ▷ initialize
26:   $m \leftarrow \text{zeros}(\text{size}(x))$  ▷ 1st moment estimate
27:   $v \leftarrow \text{zeros}(\text{size}(x))$  ▷ 2nd moment estimate
28:   $t \leftarrow 0$ 
29:  while  $t < \text{maxIter}$  do ▷ optimization loop
30:     $t \leftarrow t + 1$ 
31:     $\text{grad} \leftarrow \text{grad}_f(x)$  ▷ compute gradient
32:     $m \leftarrow \text{beta1} \times m + (1 - \text{beta1}) \times \text{grad}$  ▷ update biased 1st moment estimate
33:     $v \leftarrow \text{beta2} \times v + (1 - \text{beta2}) \times (\text{grad}^2)$  ▷ update biased 2nd moment estimate
34:     $m_{\text{hat}} \leftarrow m / (1 - \text{beta1}^t)$  ▷ correct bias in 1st moment estimate
35:     $v_{\text{hat}} \leftarrow v / (1 - \text{beta2}^t)$  ▷ correct bias in 2nd moment estimate
36:     $x \leftarrow x - \text{alpha} \times m_{\text{hat}} / (\sqrt{v_{\text{hat}}} + \text{epsilon})$  ▷ update parameters
37:  end while
38:   $f_{\text{opt}} \leftarrow f(x)$  ▷ compute optimal function value
39:  return  $x_{\text{opt}}, f_{\text{opt}}$ 
40: end procedure

```

Inputs

- f - Objective function.
- grad_f - Gradient of the objective function.
- x_0 - Initial point.
- options - Optimization options (optional).

Outputs

- x_{opt} - Optimal point.
- f_{opt} - Optimal function value.

30 AMSGrad Optimization

30.1 amsggrad Function

Description The `amsgrad` function implements the AMSGrad optimization algorithm. It performs optimization of the objective function f using the AMSGrad algorithm.

Algorithm

Algorithm 34 AMSGrad Optimization Algorithm

```

1: procedure AMSGRAD( $f$ ,  $\text{grad}_f$ ,  $x_0$ , options)
2:   if nargin < 4 then
3:     options  $\leftarrow$  struct()
4:   end if
5:   alpha  $\leftarrow$  0.01 ▷ learning rate
6:   beta1  $\leftarrow$  0.9 ▷ decay rate for 1st moment estimate
7:   beta2  $\leftarrow$  0.999 ▷ decay rate for 2nd moment estimate
8:   epsilon  $\leftarrow$   $1e - 8$  ▷ small constant for numerical stability
9:   maxIter  $\leftarrow$  1000 ▷ maximum number of iterations
10:  if isfield(options, 'alpha') then
11:    alpha  $\leftarrow$  options.alpha
12:  end if
13:  if isfield(options, 'beta1') then
14:    beta1  $\leftarrow$  options.beta1
15:  end if
16:  if isfield(options, 'beta2') then
17:    beta2  $\leftarrow$  options.beta2
18:  end if
19:  if isfield(options, 'epsilon') then
20:    epsilon  $\leftarrow$  options.epsilon
21:  end if
22:  if isfield(options, 'maxIter') then
23:    maxIter  $\leftarrow$  options.maxIter
24:  end if
25:   $x \leftarrow x_0$  ▷ initialize
26:   $m \leftarrow \text{zeros}(\text{size}(x))$  ▷ 1st moment estimate
27:   $v \leftarrow \text{zeros}(\text{size}(x))$  ▷ 2nd moment estimate
28:   $v_{\text{hat}} \leftarrow \text{zeros}(\text{size}(x))$  ▷ AMSGrad correction term
29:   $t \leftarrow 0$ 
30:  while  $t < \text{maxIter}$  do ▷ optimization loop
31:     $t \leftarrow t + 1$ 
32:     $\text{grad} \leftarrow \text{grad}_f(x)$  ▷ compute gradient
33:     $m \leftarrow \text{beta1} \times m + (1 - \text{beta1}) \times \text{grad}$  ▷ update biased 1st moment estimate
34:     $v \leftarrow \text{beta2} \times v + (1 - \text{beta2}) \times (\text{grad}^2)$  ▷ update biased 2nd moment estimate
35:     $v_{\text{hat}} \leftarrow \max(v_{\text{hat}}, v)$  ▷ update AMSGrad correction term
36:     $x \leftarrow x - \text{alpha} \times m / (\sqrt{v_{\text{hat}}} + \text{epsilon})$  ▷ update parameters
37:  end while
38:   $f_{\text{opt}} \leftarrow f(x)$  ▷ compute optimal function value
39:  return  $x_{\text{opt}}, f_{\text{opt}}$ 
40: end procedure

```

Inputs

- f - Objective function.
- grad_f - Gradient of the objective function.
- x_0 - Initial point.

- `options` - Optimization options (optional).

Outputs

- x_{opt} - Optimal point.

- f_{opt} - Optimal function value.

31 Eve Optimization

31.1 eve Function

Description The `eve` function implements the Eve optimization algorithm. It performs optimization of the objective function f using the Eve algorithm.

Algorithm

Algorithm 35 Eve Optimization Algorithm

```

1: procedure EVE( $f$ , grad_f,  $x_0$ , options)
2:   if nargin < 4 then
3:     options  $\leftarrow$  struct()
4:   end if
5:   alpha  $\leftarrow$  0.001 ▷ learning rate
6:   beta1  $\leftarrow$  0.9 ▷ decay rate for 1st moment estimate
7:   beta2  $\leftarrow$  0.999 ▷ decay rate for 2nd moment estimate
8:   epsilon  $\leftarrow$   $1e - 8$  ▷ small constant for numerical stability
9:   gamma  $\leftarrow$  0.01 ▷ smoothing parameter
10:  maxIter  $\leftarrow$  1000 ▷ maximum number of iterations
11:  if isfield(options, 'alpha') then
12:    alpha  $\leftarrow$  options.alpha
13:  end if
14:  if isfield(options, 'beta1') then
15:    beta1  $\leftarrow$  options.beta1
16:  end if
17:  if isfield(options, 'beta2') then
18:    beta2  $\leftarrow$  options.beta2
19:  end if
20:  if isfield(options, 'epsilon') then
21:    epsilon  $\leftarrow$  options.epsilon
22:  end if
23:  if isfield(options, 'gamma') then
24:    gamma  $\leftarrow$  options.gamma
25:  end if
26:  if isfield(options, 'maxIter') then
27:    maxIter  $\leftarrow$  options.maxIter
28:  end if
29:   $x \leftarrow x_0$  ▷ initialize
30:   $m \leftarrow$  zeros(size( $x$ )) ▷ 1st moment estimate
31:   $v \leftarrow$  zeros(size( $x$ )) ▷ 2nd moment estimate
32:   $t \leftarrow 0$ 
33:  while  $t <$  maxIter do ▷ optimization loop
34:     $t \leftarrow t + 1$ 
35:    grad  $\leftarrow$  grad_f( $x$ ) ▷ compute gradient
36:     $m \leftarrow$  beta1  $\times$   $m + (1 - \text{beta1}) \times$  grad ▷ update 1st moment estimate
37:     $v \leftarrow$  beta2  $\times$   $v + (1 - \text{beta2}) \times$  (grad2) ▷ update 2nd moment estimate
38:     $m_{\text{hat}} \leftarrow m / (1 - \text{beta1}^t)$  ▷ bias correction
39:     $v_{\text{hat}} \leftarrow v / (1 - \text{beta2}^t)$  ▷ bias correction
40:     $x \leftarrow x - \text{alpha} \times (m_{\text{hat}} / (\sqrt{v_{\text{hat}}} + \text{epsilon}) + \text{gamma} \times \text{grad})$  ▷ update parameters
41:  end while
42:   $f_{\text{opt}} \leftarrow f(x)$  ▷ compute optimal function value
43:  return  $x_{\text{opt}}, f_{\text{opt}}$ 
44: end procedure

```

Inputs

- f - Objective function.
- grad_f - Gradient of the objective function.
- x_0 - Initial point.
- options - Optimization options (optional).

Outputs

- x_{opt} - Optimal point.

- f_{opt} - Optimal function value.

32 Fast Adaptive Stochastic Function Accelerator (FASFA)

32.1 fasfa Function

Description The `fasfa` function implements the Fast Adaptive Stochastic Function Accelerator (FASFA) algorithm. It performs optimization of the objective function f using the FASFA algorithm.

Algorithm

Algorithm 36 FASFA Optimization Algorithm

```

1: procedure FASFA( $f$ ,  $\text{grad}_f$ ,  $x_0$ , options)
2:   if nargin < 4 then
3:     options  $\leftarrow$  struct()
4:   end if
5:   alpha  $\leftarrow$  0.001 ▷ learning rate
6:   mu  $\leftarrow$  0.8 ▷ first order momentum decay estimate
7:   nu  $\leftarrow$  0.999 ▷ second order momentum decay estimate
8:   epsilon  $\leftarrow$   $1e - 8$  ▷ small constant for numerical stability
9:   maxIter  $\leftarrow$  1000 ▷ maximum number of iterations
10:  if isfield(options, 'alpha') then
11:    alpha  $\leftarrow$  options.alpha
12:  end if
13:  if isfield(options, 'mu') then
14:    mu  $\leftarrow$  options.mu
15:  end if
16:  if isfield(options, 'nu') then
17:    nu  $\leftarrow$  options.nu
18:  end if
19:  if isfield(options, 'epsilon') then
20:    epsilon  $\leftarrow$  options.epsilon
21:  end if
22:  if isfield(options, 'maxIter') then
23:    maxIter  $\leftarrow$  options.maxIter
24:  end if
25:   $x \leftarrow x_0$  ▷ initialize
26:   $m \leftarrow \text{zeros}(\text{size}(x))$  ▷ first moment vector
27:   $n \leftarrow \text{zeros}(\text{size}(x))$  ▷ second moment vector
28:   $t \leftarrow 0$  ▷ timestep/iteration
29:  while  $t < \text{maxIter}$  do ▷ optimization loop
30:     $t \leftarrow t + 1$ 
31:     $\text{grad} \leftarrow \text{grad}_f(x)$  ▷ compute gradient
32:     $m \leftarrow \text{mu} \times m + (1 - \text{mu}) \times \text{grad}$  ▷ update biased first moment estimation
33:     $n \leftarrow \text{nu} \times n + (1 - \text{nu}) \times (\text{grad}^2)$  ▷ update biased second moment estimation
34:     $m_{\text{hat}} \leftarrow m / (1 - \text{mu}^t)$  ▷ compute raw moment estimates
35:     $n_{\text{hat}} \leftarrow n / (1 - \text{nu}^t)$ 
36:     $x \leftarrow x - \text{alpha} \times m_{\text{hat}} / (\sqrt{n_{\text{hat}}} + \text{epsilon})$  ▷ implement FASFA update rule
37:  end while
38:   $f_{\text{opt}} \leftarrow f(x)$  ▷ compute optimal function value
39:  return  $x_{\text{opt}}$ ,  $f_{\text{opt}}$ 
40: end procedure

```

Inputs

- f - Objective function.

- `grad_f` - Gradient of the objective function.
- `x0` - Initial point.
- `options` - Optimization options (optional).

Outputs

- x_{opt} - Optimal point.
- f_{opt} - Optimal function value.

33 Gradient Descent with Momentum

33.1 `gradientDescentWithMomentum` Function

Description The `gradientDescentWithMomentum` function implements the Gradient Descent with Momentum optimization algorithm. It performs optimization of the objective function f using the Gradient Descent with Momentum algorithm.

Algorithm

Algorithm 37 Gradient Descent with Momentum Optimization Algorithm

```

1: procedure GRADIENTDESCENTWITHMOMENTUM( $f$ ,  $\text{grad}_f$ ,  $x_0$ ,  $\text{options}$ )
2:   if nargin < 4 then
3:      $\text{options} \leftarrow \text{struct}()$ 
4:   end if
5:    $\alpha \leftarrow 0.01$  ▷ learning rate
6:    $\beta \leftarrow 0.9$  ▷ momentum coefficient
7:    $\text{maxIter} \leftarrow 1000$  ▷ maximum number of iterations
8:   if isfield( $\text{options}$ , 'alpha') then
9:      $\alpha \leftarrow \text{options.alpha}$ 
10:  end if
11:  if isfield( $\text{options}$ , 'beta') then
12:     $\beta \leftarrow \text{options.beta}$ 
13:  end if
14:  if isfield( $\text{options}$ , 'maxIter') then
15:     $\text{maxIter} \leftarrow \text{options.maxIter}$ 
16:  end if
17:   $x \leftarrow x_0$  ▷ initialize
18:   $v \leftarrow \text{zeros}(\text{size}(x))$  ▷ momentum
19:  for  $t = 1$  to  $\text{maxIter}$  do ▷ optimization loop
20:     $\text{grad} \leftarrow \text{grad}_f(x)$  ▷ compute gradient
21:     $v \leftarrow \beta \times v + \alpha \times \text{grad}$  ▷ update momentum
22:     $x \leftarrow x - v$  ▷ update parameters
23:  end for
24:   $f_{\text{opt}} \leftarrow f(x)$  ▷ compute optimal function value
25:  return  $x_{\text{opt}}, f_{\text{opt}}$ 
26: end procedure

```

Inputs

- f - Objective function.
- `grad_f` - Gradient of the objective function.
- `x0` - Initial point.
- `options` - Optimization options (optional).

Outputs

- x_{opt} - Optimal point.
- f_{opt} - Optimal function value.

34 Mini-Batch Gradient Descent with Momentum

34.1 miniBatchGradientDescentWithMomentum Function

Description The `miniBatchGradientDescentWithMomentum` function implements the Mini-Batch Gradient Descent with Momentum optimization algorithm. It performs optimization of the objective function f using the Mini-Batch Gradient Descent with Momentum algorithm.

Algorithm

Algorithm 38 Mini-Batch Gradient Descent with Momentum Optimization Algorithm

```

1: procedure MINIBATCHGRADIENTDESCENTWITHMOMENTUM( $f$ ,  $\text{grad}_f$ ,  $x_0$ , options)
2:   if nargin < 4 then
3:     options  $\leftarrow$  struct()
4:   end if
5:   alpha  $\leftarrow$  0.01 ▷ learning rate
6:   beta  $\leftarrow$  0.9 ▷ momentum coefficient
7:   batchSize  $\leftarrow$  32 ▷ batch size
8:   maxIter  $\leftarrow$  1000 ▷ maximum number of iterations
9:   if isfield(options, 'alpha') then
10:    alpha  $\leftarrow$  options.alpha
11:  end if
12:  if isfield(options, 'beta') then
13:    beta  $\leftarrow$  options.beta
14:  end if
15:  if isfield(options, 'batchSize') then
16:    batchSize  $\leftarrow$  options.batchSize
17:  end if
18:  if isfield(options, 'maxIter') then
19:    maxIter  $\leftarrow$  options.maxIter
20:  end if
21:   $x \leftarrow x_0$  ▷ initialize
22:   $v \leftarrow \text{zeros}(\text{size}(x))$  ▷ momentum
23:  for  $t = 1$  to maxIter do ▷ optimization loop
24:    rand_indices  $\leftarrow$  randperm(length( $x_0$ ), batchSize) ▷ sample mini-batch
25:    grad_sum  $\leftarrow$  zeros(size( $x$ ))
26:    for  $i = 1$  to length(rand_indices) do grad  $\leftarrow$  grad_f( $x(:, \text{rand\_indices}(i))$ )
27:      grad_sum  $\leftarrow$  grad_sum + grad
28:    end for
29:
30:    grad_avg  $\leftarrow$  grad_sum/batchSize ▷ update momentum
31:     $v \leftarrow \text{beta} \times v + \text{alpha} \times \text{grad\_avg}$ 
32:     $x \leftarrow x - v$  ▷ update parameters
33:
34:     $f_{\text{opt}} \leftarrow f(x)$  ▷ compute optimal function value
35:  return  $x_{\text{opt}}, f_{\text{opt}}$ 
36:  =0

```

Inputs

- f - Objective function.
- grad_f - Gradient of the objective function.

- x_0 - Initial point.
- `options` - Optimization options (optional).

Outputs

- x_{opt} - Optimal point.
- f_{opt} - Optimal function value.

35 Nesterov Momentum

35.1 nestrovMomentum Function

Description The `nestrovMomentum` function implements the Nesterov Momentum optimization algorithm. It performs optimization of the objective function f using the Nesterov Momentum algorithm.

Algorithm

Algorithm 39 Nesterov Momentum Optimization Algorithm

```

1: procedure NESTROVMOMENTUM( $f$ , grad_f,  $x_0$ , options)
2:   if nargin < 4 then
3:     options ← struct()
4:   end if
5:   alpha ← 0.01                                     ▷ learning rate
6:   beta ← 0.9                                       ▷ momentum coefficient
7:   maxIter ← 1000                                   ▷ maximum number of iterations
8:   if isfield(options, 'alpha') then
9:     alpha ← options.alpha
10:  end if
11:  if isfield(options, 'beta') then
12:    beta ← options.beta
13:  end if
14:  if isfield(options, 'maxIter') then
15:    maxIter ← options.maxIter
16:  end if
17:   $x \leftarrow x_0$                                   ▷ initialize
18:   $v \leftarrow \text{zeros}(\text{size}(x))$                   ▷ velocity
19:  for  $t = 1$  to maxIter do                          ▷ optimization loop
20:    lookahead_x ←  $x - \beta \times v$                 ▷ compute gradient at lookahead point
21:    grad ← grad_f(lookahead_x)
22:     $v \leftarrow \beta \times v + \alpha \times \text{grad}$     ▷ update velocity
23:     $x \leftarrow x - v$                               ▷ update parameters
24:  end for
25:   $f_{\text{opt}} \leftarrow f(x)$                           ▷ compute optimal function value
26:  return  $x_{\text{opt}}, f_{\text{opt}}$ 
27: end procedure

```

Inputs

- f - Objective function.
- `grad_f` - Gradient of the objective function.
- x_0 - Initial point.
- `options` - Optimization options (optional).

Outputs

- x_{opt} - Optimal point.
- f_{opt} - Optimal function value.

36 RMSProp

36.1 rmsprop Function

Description The `rmsprop` function implements the RMSProp optimization algorithm. It performs optimization of the objective function f using the RMSProp algorithm.

Algorithm

Algorithm 40 RMSProp Optimization Algorithm

```

1: procedure RMSPROP( $f$ , grad_f,  $x_0$ , options)
2:   if nargin < 4 then
3:     options ← struct()
4:   end if
5:   alpha ← 0.001 ▷ learning rate
6:   beta ← 0.9 ▷ decay rate
7:   epsilon ← 1e - 8 ▷ small constant for numerical stability
8:   maxIter ← 1000 ▷ maximum number of iterations
9:   if isfield(options, 'alpha') then
10:    alpha ← options.alpha
11:  end if
12:  if isfield(options, 'beta') then
13:    beta ← options.beta
14:  end if
15:  if isfield(options, 'epsilon') then
16:    epsilon ← options.epsilon
17:  end if
18:  if isfield(options, 'maxIter') then
19:    maxIter ← options.maxIter
20:  end if
21:   $x \leftarrow x_0$  ▷ initialize
22:  E_g ← zeros(size( $x$ )) ▷ running average of squared gradients
23:  for  $t = 1$  to maxIter do ▷ optimization loop
24:    grad ← grad_f( $x$ ) ▷ compute gradient
25:    E_g ← beta × E_g + (1 - beta) × grad2 ▷ compute squared gradient
26:     $x \leftarrow x - \text{alpha} \times \text{grad} / (\sqrt{\text{E\_g}} + \text{epsilon})$  ▷ update parameters
27:  end for
28:   $f_{\text{opt}} \leftarrow f(x)$  ▷ compute optimal function value
29:  return  $x_{\text{opt}}, f_{\text{opt}}$ 
30: end procedure

```

Inputs

- f - Objective function.
- grad_f - Gradient of the objective function.
- x_0 - Initial point.
- options - Optimization options (optional).

Outputs

- x_{opt} - Optimal point.
- f_{opt} - Optimal function value.

37 Stochastic Gradient Descent with Momentum

37.1 stochasticGradientDescentWithMomentum Function

Description The `stochasticGradientDescentWithMomentum` function implements the Stochastic Gradient Descent with Momentum optimization algorithm. It performs optimization of the objective function f using the Stochastic Gradient Descent with Momentum algorithm.

Algorithm

Algorithm 41 Stochastic Gradient Descent with Momentum Optimization Algorithm

```

1: procedure STOCHASTICGRADIENTDESCENTWITHMOMENTUM( $f$ ,  $\text{grad}_f$ ,  $x_0$ , options)
2:   if nargin < 4 then
3:     options  $\leftarrow$  struct()
4:   end if
5:   alpha  $\leftarrow$  0.01 ▷ learning rate
6:   beta  $\leftarrow$  0.9 ▷ momentum coefficient
7:   maxIter  $\leftarrow$  1000 ▷ maximum number of iterations
8:   if isfield(options, 'alpha') then
9:     alpha  $\leftarrow$  options.alpha
10:  end if
11:  if isfield(options, 'beta') then
12:    beta  $\leftarrow$  options.beta
13:  end if
14:  if isfield(options, 'maxIter') then
15:    maxIter  $\leftarrow$  options.maxIter
16:  end if
17:   $x \leftarrow x_0$  ▷ initialize
18:   $v \leftarrow \text{zeros}(\text{size}(x))$  ▷ momentum
19:  for  $t = 1$  to maxIter do ▷ optimization loop
20:    rand_index  $\leftarrow$  randi([1, length( $x_0$ )]) ▷ sample a random data point
21:    grad  $\leftarrow$  grad_f( $x(:, \text{rand\_index})$ ) ▷ compute gradient
22:     $v \leftarrow \text{beta} \times v + \text{alpha} \times \text{grad}$  ▷ update momentum
23:     $x \leftarrow x - v$  ▷ update parameters
24:  end for
25:   $f_{\text{opt}} \leftarrow f(x)$  ▷ compute optimal function value
26:  return  $x_{\text{opt}}, f_{\text{opt}}$ 
27: end procedure

```

Inputs

- f - Objective function.
- grad_f - Gradient of the objective function.
- x_0 - Initial point.
- options - Optimization options (optional).

Outputs

- x_{opt} - Optimal point.
- f_{opt} - Optimal function value.

38 Exponentially Weighted Online Convex Optimization

38.1 expWeightedOCO Function

Description The `expWeightedOCO` function performs exponentially weighted online convex optimization. It optimizes the objective function f with its gradient `gradf` starting from the initial guess x_0 using the stepsize α for a maximum of `maxIter` iterations.

Algorithm

Algorithm 42 Exponentially Weighted Online Convex Optimization

```

1: procedure EXPWEIGHTEDOCO( $f$ , gradf,  $x_0$ ,  $\alpha$ , maxIter)
2:    $x \leftarrow x_0$  ▷ Initialize parameters
3:   history  $\leftarrow$  zeros(maxIter, 1) ▷ Initialize history of function values
4:   weights  $\leftarrow$  ones(maxIter, 1) ▷ Initialize weights
5:   for  $k = 1$  to maxIter do
6:      $f_k \leftarrow f(k, x)$  ▷ Current function value
7:      $\text{grad}f_k \leftarrow \text{grad}f(k, x)$  ▷ Current gradient
8:      $\text{weights}(k) \leftarrow \exp(-\alpha \cdot f_k)$  ▷ Update weight
9:      $x \leftarrow x - \alpha \cdot \text{grad}f_k \cdot \text{weights}(k)$  ▷ Update parameters
10:    history( $k$ )  $\leftarrow f_k$  ▷ Store current function value in history
11:  end for
12:  return  $x$ , history
13: end procedure

```

Inputs

- f - Function handle of the objective function (sequence of functions).
- `gradf` - Function handle of the gradient of f (sequence of gradients).
- x_0 - Initial guess.
- α - Step size.
- `maxIter` - Maximum number of iterations.

Outputs

- x - Optimized parameters.
- `history` - History of function values.

38.2 onlineGradientDescent Function

Description The `onlineGradientDescent` function performs online gradient descent optimization. It optimizes the objective function f with its gradient `gradf` starting from the initial guess x_0 using the stepsize α for a maximum of `maxIter` iterations.

Algorithm

Algorithm 43 Online Gradient Descent

```

1: procedure ONLINEGRADIENTDESCENT( $f$ , gradf,  $x_0$ ,  $\alpha$ , maxIter)
2:    $x \leftarrow x_0$                                      ▷ Initialize parameters
3:   history  $\leftarrow$  zeros(maxIter, 1)                ▷ Initialize history of function values
4:   for  $k = 1$  to maxIter do
5:      $f_k \leftarrow f(k, x)$                           ▷ Current function value
6:     gradf $_k \leftarrow$  gradf( $k, x$ )                  ▷ Current gradient
7:      $x \leftarrow x - \alpha \cdot$  gradf $_k$              ▷ Update parameters
8:     history( $k$ )  $\leftarrow$   $f_k$                        ▷ Store current function value in history
9:   end for
10:  return  $x$ , history
11: end procedure

```

Inputs

- f - Function handle of the objective function (sequence of functions).
- gradf - Function handle of the gradient of f (sequence of gradients).
- x_0 - Initial guess.
- α - Step size.
- maxIter - Maximum number of iterations.

Outputs

- x - Optimized parameters.
- history - History of function values.

39 Online Newton Step

39.1 onlineNewtonStep Function

Description The `onlineNewtonStep` function performs online Newton step optimization. It optimizes the objective function f with its gradient `gradf` and Hessian `hessf` starting from the initial guess x_0 using the stepsize α for a maximum of `maxIter` iterations.

Algorithm

Algorithm 44 Online Newton Step

```

1: procedure ONLINENEWTONSTEP( $f$ , gradf, hessf,  $x_0$ ,  $\alpha$ , maxIter)
2:    $x \leftarrow x_0$                                      ▷ Initialize parameters
3:   history  $\leftarrow$  zeros(maxIter, 1)                ▷ Initialize history of function values
4:   for  $k = 1$  to maxIter do
5:      $f_k \leftarrow f(k, x)$                           ▷ Current function value
6:     gradf $_k \leftarrow$  gradf( $k, x$ )                  ▷ Current gradient
7:     hessf $_k \leftarrow$  hessf( $k, x$ )                  ▷ Current Hessian
8:      $x \leftarrow x - \alpha \cdot$  (hessf $_k \setminus$  gradf $_k$ )  ▷ Update parameters
9:     history( $k$ )  $\leftarrow$   $f_k$                        ▷ Store current function value in history
10:  end for
11:  return  $x$ , history
12: end procedure

```

Inputs

- f - Function handle of the objective function (sequence of functions).

- `gradf` - Function handle of the gradient of f (sequence of gradients).
- `hessf` - Function handle of the Hessian of f (sequence of Hessians).
- x_0 - Initial guess.
- α - Step size.
- `maxIter` - Maximum number of iterations.

Outputs

- x - Optimized parameters.
- `history` - History of function values.

40 Stochastic Gradient Descent

40.1 `stochasticGradientDescent` Function

Description The `stochasticGradientDescent` function performs stochastic gradient descent optimization. It optimizes the objective function f with its gradient `gradf` starting from the initial guess x_0 using the stepsize α for a maximum of `maxIter` iterations and batch size `batchSize`.

Algorithm

Algorithm 45 Stochastic Gradient Descent

```

1: procedure STOCHASTICGRADIENTDESCENT( $f$ , gradf,  $x_0$ ,  $\alpha$ , maxIter, batchSize)
2:    $x \leftarrow x_0$  ▷ Initialize parameters
3:   history  $\leftarrow$  zeros(maxIter, 1) ▷ Initialize history of function values
4:    $n \leftarrow$  numel( $f$ ) ▷ Total number of functions
5:   for  $k = 1$  to maxIter do
6:      $\text{idx} \leftarrow$  randperm( $n$ , batchSize) ▷ Randomly sample batch indices
7:     gradBatch  $\leftarrow$  zeros(size( $x$ )) ▷ Initialize gradient batch
8:     for  $j = 1$  to batchSize do
9:       gradBatch  $\leftarrow$  gradBatch + gradf(idx( $j$ ),  $x$ ) ▷ Accumulate gradients
10:    end for
11:    gradBatch  $\leftarrow$  gradBatch/batchSize ▷ Average gradient
12:     $x \leftarrow x - \alpha \cdot$  gradBatch ▷ Update parameters
13:    history( $k$ )  $\leftarrow$  mean(arrayfun(@(i)f( $i$ ,  $x$ ), idx)) ▷ Store current function value in history
14:  end for
15:  return  $x$ , history
16: end procedure

```

Inputs

- f - Function handle of the objective function (sequence of functions).
- `gradf` - Function handle of the gradient of f (sequence of gradients).
- x_0 - Initial guess.
- α - Step size.
- `maxIter` - Maximum number of iterations.
- `batchSize` - Size of the batch for stochastic updates.

Outputs

- x - Optimized parameters.
- `history` - History of function values.

41 Conditional Gradient

41.1 conditionalGradient Function

Description The `conditionalGradient` function implements the conditional gradient (Frank-Wolfe) method. It optimizes the objective function f with its gradient `gradf` starting from the initial guess x_0 using linear constraints defined by A and b for a maximum of `maxIter` iterations with a tolerance of `tol`.

Algorithm

Algorithm 46 Conditional Gradient

```

1: procedure CONDITIONALGRADIENT( $f$ , gradf,  $x_0$ ,  $A$ ,  $b$ , maxIter, tol)
2:    $x \leftarrow x_0$  ▷ Initialize parameters
3:   history  $\leftarrow$  zeros(maxIter, 1) ▷ Initialize history of function values
4:   for  $k = 1$  to maxIter do
5:     gradient  $\leftarrow$  gradf( $x$ ) ▷ Compute gradient
6:      $s \leftarrow$  linearOracle(gradient,  $A$ ,  $b$ ) ▷ Solve linear problem
7:      $\alpha \leftarrow \frac{2}{k+2}$  ▷ Line search stepsize
8:      $x \leftarrow (1 - \alpha) \cdot x + \alpha \cdot s$  ▷ Update parameters
9:     history( $k$ )  $\leftarrow$   $f(x)$  ▷ Store current function value in history
10:    if  $\|\text{gradf}(x)\| < \text{tol}$  then ▷ Check for convergence
11:      history  $\leftarrow$  history(1 :  $k$ ) ▷ Trim history
12:    break
13:  end if
14: end for
15: return  $x$ , history
16: end procedure

```

Inputs

- f - Function handle of the objective function.
- `gradf` - Function handle of the gradient of f .
- x_0 - Initial guess.
- A - Matrix defining the linear constraints.
- b - Vector defining the linear constraints.
- `maxIter` - Maximum number of iterations.
- `tol` - Tolerance for convergence.

Outputs

- x - Optimized parameters.
- `history` - History of function values.

41.2 linearOracle Function

Description The `linearOracle` function solves the linear problem to find the update direction in the conditional gradient method.

Algorithm The linear problem is solved using the `linprog` function, which is a linear programming solver.

Inputs

- `gradient` - Gradient vector.
- A - Matrix defining the linear constraints.
- b - Vector defining the linear constraints.

41.2.1 Output

- s - Update direction.

42 Conditional Gradient Projection-Free

42.1 conditionalGradientProjectionFree Function

Description The `conditionalGradientProjectionFree` function implements the projection-free conditional gradient method. It optimizes the objective function f with its gradient `gradf` starting from the initial guess x_0 using linear constraints defined by A and b for a maximum of `maxIter` iterations with a tolerance of `tol`.

Algorithm

Algorithm 47 Conditional Gradient Projection-Free

```

1: procedure CONDITIONALGRADIENTPROJECTIONFREE( $f$ ,  $\text{gradf}$ ,  $x_0$ ,  $A$ ,  $b$ ,  $\text{maxIter}$ ,  $\text{tol}$ )
2:    $x \leftarrow x_0$                                      ▷ Initialize parameters
3:    $\text{history} \leftarrow \text{zeros}(\text{maxIter}, 1)$            ▷ Initialize history of function values
4:   for  $k = 1$  to  $\text{maxIter}$  do
5:      $\text{gradient} \leftarrow \text{gradf}(x)$                  ▷ Compute gradient
6:      $s \leftarrow \text{linearOracle}(\text{gradient}, A, b)$      ▷ Solve linear problem
7:      $\alpha \leftarrow \frac{2}{k+2}$                          ▷ Line search stepsize
8:      $x \leftarrow (1 - \alpha) \cdot x + \alpha \cdot s$    ▷ Update parameters
9:      $\text{history}(k) \leftarrow f(x)$                      ▷ Store current function value in history
10:    if  $\|\text{gradf}(x)\| < \text{tol}$  then                 ▷ Check for convergence
11:       $\text{history} \leftarrow \text{history}(1 : k)$            ▷ Trim history
12:    break
13:  end if
14: end for
15:  return  $x$ ,  $\text{history}$ 
16: end procedure

```

Inputs

- f - Function handle of the objective function.
- `gradf` - Function handle of the gradient of f .
- x_0 - Initial guess.
- A - Matrix defining the linear constraints.
- b - Vector defining the linear constraints.
- `maxIter` - Maximum number of iterations.
- `tol` - Tolerance for convergence.

Outputs

- x - Optimized parameters.
- `history` - History of function values.

43 Linear Oracle

43.1 linearOracle Function

Description The `linearOracle` function solves the linear problem to find the update direction subject to linear constraints.

43.1.1 Function

Algorithm 48 Linear Oracle

```

1: procedure LINEARORACLE(gradient, A, b)
2:    $K \leftarrow \text{length}(b)$ 
3:    $f \leftarrow \text{gradient}$ 
4:    $\text{options} \leftarrow \text{optimoptions}('linprog', 'Display', 'none')$ 
5:    $s \leftarrow \text{linprog}(f, A, b, [], [], \text{zeros}(K, 1), [], \text{options})$ 
6:   return  $s$ 
7: end procedure

```

▷ Number of constraints
 ▷ Objective function
 ▷ Options for solver
 ▷ Solve linear problem
 ▷ Return update direction

Inputs

- *gradient* - Gradient vector.
- *A* - Matrix defining the linear constraints.
- *b* - Vector defining the linear constraints.

43.1.2 Output

- *s* - Update direction.

44 Online Conditional Gradient

44.1 onlineConditionalGradient Function

Description The `onlineConditionalGradient` function implements the online conditional gradient method to optimize a function subject to linear constraints.

44.1.1 Function

Algorithm 49 Online Conditional Gradient

```

1: procedure ONLINECONDITIONALGRADIENT(f, gradf, x0, A, b, maxIter, tol,  $\eta$ )
2:    $x \leftarrow x0$ 
3:    $\text{history} \leftarrow \text{zeros}(\text{maxIter}, 1)$ 
4:   for  $k = 1$  to  $\text{maxIter}$  do
5:      $\text{gradient} \leftarrow \text{gradf}(x)$ 
6:      $s \leftarrow \text{linearOracle}(\text{gradient}, A, b)$ 
7:      $x \leftarrow (1 - \eta) \cdot x + \eta \cdot s$ 
8:      $\text{history}(k) \leftarrow f(x)$ 
9:     if  $\|\text{gradf}(x)\| < \text{tol}$  then
10:       $\text{history} \leftarrow \text{history}(1 : k)$ 
11:      break
12:     end if
13:   end for
14:   return  $x$ ,  $\text{history}$ 
15: end procedure

```

▷ Initial guess
 ▷ History of function values
 ▷ Compute gradient
 ▷ Solve linear problem
 ▷ Update parameters
 ▷ Record function value
 ▷ Check convergence
 ▷ Return optimized parameters and history

Inputs

- *f* - Function handle of the objective function.
- *gradf* - Function handle of the gradient of the objective function.
- *x0* - Initial guess.
- *A* - Matrix defining the linear constraints.

- b - Vector defining the linear constraints.
- `maxIter` - Maximum number of iterations.
- `tol` - Tolerance for convergence.
- η - Step size.

Outputs

- x - Optimized parameters.
- `history` - History of function values.

45 Adaptive Gradient Descent

45.1 `adaptiveGradientDescent` Function

Description The `adaptiveGradientDescent` function performs adaptive gradient descent optimization by adjusting the stepsize based on the adaptation parameter.

45.1.1 Function

Algorithm 50 Adaptive Gradient Descent

```

1: procedure ADAPTIVEGRADIENTDESCENT( $f$ ,  $\text{grad}f$ ,  $x_0$ ,  $\alpha$ ,  $\beta$ ,  $\text{maxIter}$ )
2:    $x \leftarrow x_0$  ▷ Initial guess
3:    $\text{history} \leftarrow \text{zeros}(\text{maxIter}, 1)$  ▷ History of function values
4:    $\text{gradHist} \leftarrow \text{zeros}(\text{size}(x_0))$  ▷ History of gradients
5:   for  $k = 1$  to  $\text{maxIter}$  do
6:      $f_k \leftarrow f(k, x)$  ▷ Current function
7:      $\text{grad}fk \leftarrow \text{grad}f(k, x)$  ▷ Current gradient
8:      $\text{gradHist} \leftarrow \text{gradHist} + \text{grad}fk^2$  ▷ Update gradient history
9:      $\text{adjAlpha} \leftarrow \alpha / (\sqrt{\text{gradHist}} + \beta)$  ▷ Adjust stepsize
10:     $x \leftarrow x - \text{adjAlpha} \cdot \text{grad}fk$  ▷ Update parameters
11:     $\text{history}(k) \leftarrow f_k$  ▷ Record function value
12:  end for
13:  return  $x$ ,  $\text{history}$  ▷ Return optimized parameters and history
14: end procedure

```

Inputs

- f - Function handle of the objective function.
- $\text{grad}f$ - Function handle of the gradient of the objective function.
- x_0 - Initial guess.
- α - Initial stepsize.
- β - Adaptation parameter.
- `maxIter` - Maximum number of iterations.

46 Follow The Perturbed Leader

46.1 `followPerturbedLeader` Function

Description The `followPerturbedLeader` function performs Follow-The-Perturbed-Leader optimization by updating the parameters based on the gradient with added perturbations.

46.1.1 Function

Algorithm 51 Follow The Perturbed Leader

```

1: procedure FOLLOWPERTURBEDLEADER( $f$ ,  $\text{gradf}$ ,  $x_0$ ,  $\alpha$ ,  $\sigma$ ,  $\text{maxIter}$ )
2:    $x \leftarrow x_0$  ▷ Initial guess
3:    $\text{history} \leftarrow \text{zeros}(\text{maxIter}, 1)$  ▷ History of function values
4:    $G \leftarrow \text{zeros}(\text{length}(x_0), \text{maxIter})$  ▷ Matrix to store gradients
5:   for  $k = 1$  to  $\text{maxIter}$  do
6:      $\text{perturbation} \leftarrow \sigma \cdot \text{randn}(\text{size}(x_0))$  ▷ Add perturbation
7:      $G(:, k) \leftarrow \text{gradf}(k, x) + \text{perturbation}$  ▷ Update gradients with perturbation
8:      $\text{sumG} \leftarrow \text{sum}(G(:, 1 : k), 2)$  ▷ Compute sum of gradients
9:      $x \leftarrow -(\alpha/k) \cdot \text{sumG}$  ▷ Update parameters
10:     $\text{history}(k) \leftarrow f(k, x)$  ▷ Record function value
11:  end for
12:  return  $x$ ,  $\text{history}$  ▷ Return optimized parameters and history
13: end procedure

```

Inputs

- f - Function handle of the objective function.
- gradf - Function handle of the gradient of the objective function.
- x_0 - Initial guess.
- α - Step size.
- σ - Perturbation parameter.
- maxIter - Maximum number of iterations.

Outputs

- x - Optimized parameters.
- history - History of function values.

47 Online Mirror Descent

47.1 onlineMirrorDescent Function

Description The `onlineMirrorDescent` function performs online mirror descent optimization by updating the parameters using a mirror map and its inverse.

47.1.1 Function

Algorithm 52 Online Mirror Descent

```

1: procedure ONLINEMIRRORDESCENT( $f$ ,  $\text{gradf}$ ,  $x_0$ ,  $\alpha$ ,  $\text{maxIter}$ ,  $\text{mirrorMap}$ ,  $\text{invMirrorMap}$ )
2:    $x \leftarrow x_0$  ▷ Initial guess
3:    $\text{history} \leftarrow \text{zeros}(\text{maxIter}, 1)$  ▷ History of function values
4:   for  $k = 1$  to  $\text{maxIter}$  do
5:      $fk \leftarrow f(k, x)$  ▷ Current function
6:      $\text{gradfk} \leftarrow \text{gradf}(k, x)$  ▷ Current gradient
7:      $z \leftarrow \text{mirrorMap}(x) - \alpha \cdot \text{gradfk}$  ▷ Mirror descent step
8:      $x \leftarrow \text{invMirrorMap}(z)$  ▷ Update parameters
9:      $\text{history}(k) \leftarrow fk$  ▷ Record function value
10:  end for
11:  return  $x$ ,  $\text{history}$  ▷ Return optimized parameters and history
12: end procedure

```

Inputs

- f - Function handle of the objective function.
- gradf - Function handle of the gradient of the objective function.
- x_0 - Initial guess.
- α - Step size.
- maxIter - Maximum number of iterations.
- mirrorMap - Mirror map function handle.
- invMirrorMap - Inverse mirror map function handle.

Outputs

- x - Optimized parameters.
- history - History of function values.

48 Regularized Follow-The-Leader (RFTL) Optimization

48.1 `rftlAlgorithm` Function

Description The `rftlAlgorithm` function performs Regularized Follow-The-Leader (RFTL) optimization by updating the parameters using a stepsize and regularization parameter.

48.1.1 Function

Algorithm 53 Regularized Follow-The-Leader (RFTL) Optimization

```

1: procedure RFTLALGORITHM( $f, \text{gradf}, x_0, \alpha, \lambda, \text{maxIter}$ )
2:    $x \leftarrow x_0$  ▷ Initial guess
3:    $\text{history} \leftarrow \text{zeros}(\text{maxIter}, 1)$  ▷ History of function values
4:    $G \leftarrow \text{zeros}(\text{length}(x_0), \text{maxIter})$  ▷ Matrix to store gradients
5:   for  $k = 1$  to  $\text{maxIter}$  do
6:      $G(:, k) \leftarrow \text{gradf}(k, x)$  ▷ Compute gradient
7:      $\text{sumG} \leftarrow \text{sum}(G(:, 1 : k), 2)$  ▷ Compute cumulative sum of gradients
8:      $x \leftarrow -\frac{\alpha}{k+\lambda} \cdot \text{sumG}$  ▷ Update parameters
9:      $\text{history}(k) \leftarrow f(k, x)$  ▷ Record function value
10:  end for
11:  return  $x, \text{history}$  ▷ Return optimized parameters and history
12: end procedure

```

Inputs

- f - Function handle of the objective function.
- gradf - Function handle of the gradient of the objective function.
- x_0 - Initial guess.
- α - Step size.
- λ - Regularization parameter.
- maxIter - Maximum number of iterations.

Outputs

- x - Optimized parameters.
- history - History of function values.

49 Discussion

Optimization has become an versatile tool across multiple disciplines, including deep learning to finance and physics. As the complexity of problems in these domains continues to grow, the need for efficient and well-organized optimization toolboxes becomes increasingly needed [4].

The OCO is a range of algorithms made specifically for solving these optimization problems in dynamic, data-driven environments where decisions must be made sequentially without knowledge of future data. This is different from a purely stochastic process, which is completely memoryless. Among these OCO algorithms, Stochastic Gradient Descent (SGD) stands out as a workhorse, widely employed in machine learning, particularly with vast datasets or when the objective function consists of numerous individual loss functions. The efficiency stems from updating parameters based on gradients computed from mini-batches of data, making it well-suited for online learning scenarios [8, 4].

ExpWeightedOCO, on the other hand, is adept at handling changing data distributions or non-stationary environments. By assigning exponentially decreasing weights to historical data, it adapts to evolving conditions by prioritizing recent observations. This makes it valuable for scenarios where data patterns fluctuate over time [4].

Meanwhile, Online Gradient Descent (OGD) offers a fundamental approach in OCO, providing regret guarantees in convex online learning settings. OGD updates parameters iteratively by moving them in the direction of the negative gradient of the current loss function. It's a simple yet powerful algorithm applicable across a wide spectrum of problems, including online prediction and dynamic resource allocation [4].

For problems with strongly convex objectives and Lipschitz continuous Hessians, the Online Newton Step tends to be better. Leveraging gradient and curvature information, it achieves faster convergence compared to gradient-based methods, making it suitable for scenarios where such properties hold [4].

This toolbox further includes more modern solvers, such as Adam and AMSGrad. The former combines properties of AdaGrad and RMSProp, which are also in this toolbox, in order to achieve fast convergence. The latter was created for one dimensional and strongly convex settings, where Adam struggles to converge [6].

These algorithms collectively form a versatile toolbox for addressing a broad array of optimization challenges in online learning, ranging from simple convex problems to more intricate and structured optimization tasks. Their selection hinges on factors like the convexity of the objective function, the availability of gradient or Hessian information, and the computational resources at hand, ensuring tailored solutions to diverse optimization scenarios. However, there can be other tools built into the toolbox that can help with optimization and basic machine learning problems. [7]

There are many additions to this toolbox that are in the works. Hyperparameterization and learning rate scheduling are both techniques in machine learning and optimization that aimed at improving model performance and convergence. While they share the overarching goal of optimizing model parameters, they operate at different levels and address distinct aspects of the optimization process. In essence, they try to accomplish the same goal in two different ways [4, 2].

Hyperparameterization involves tuning the hyperparameters of a learning algorithm, such as regularization strength or kernel parameters in a support vector machine. One of the most effective use cases for this approach is with the deep convolutional neural networks, and furthermore on generative adversarial networks. Usually, this is an automated process aimed at finding the optimal configuration of hyperparameters that leads to improved model performance, generalization, and convergence. Hyperparameterization is crucial for fine-tuning models to specific datasets and problem domains, as it directly influences the model's capacity to learn and generalize from data [9].

Learning rate scheduling, on the other hand, focuses on adapting the learning rate during training to facilitate efficient optimization and convergence. It is the simpler approach of the two, and can have a low computational cost. The learning rate determines the stepsize taken in the parameter space during some gradient-based optimization algorithms. Learning rate scheduling techniques adjust the learning rate over time, often based on heuristics, to balance between rapid convergence in the early stages of training and fine-tuning towards the end. Common scheduling strategies include exponential decay, step decay, polynomial decay, cosine annealing, and adaptive methods like Adam and RMSprop. The most effective versions of this technique often involve geometrically annealing a repeating pattern [2].

While hyperparameterization and learning rate scheduling share the goal of improving optimization performance, they operate at different levels of abstraction and address different aspects of the optimization process. Hyperparameterization involves tuning model-specific parameters that influence the learning process [9], while learning rate scheduling focuses on dynamically adjusting the optimization algorithm's behavior during training [2]. Currently, both of these are being added to the toolbox.

In some scenarios, one of these can be better than the other, but that is largely depending on the nature of the problem and available resources. For instance, in settings with limited computational resources or when manual intervention is feasible, hyperparameterization may offer more flexibility and control over model performance. On the other hand, learning rate scheduling can be particularly beneficial in deep learning applications, where finding the optimal learning rate manually can be challenging due to the complexity of the model and the large-scale datasets involved. Adaptive methods like Adam and RMSprop automatically adjust the learning rate based on the gradient’s statistics anyways, alleviating the need for extensive hyperparameter tuning [2, 9].

Approximately five years ago, Jagdeep Bhatia proposed new algorithms for machine learning. They use a principle called interactive machine learning using randomized counter examples. Among these algorithms, arbitrary learning stands out for its flexibility, allowing users to provide diverse forms of feedback or labels to guide the learning process. This approach is particularly beneficial in scenarios where traditional labeled datasets are scarce or expensive to obtain, enabling the model to learn from various sources of information, including domain knowledge and user preferences. Though these algorithms are somewhat outside the domain of the OCO, perhaps they could be used in conjunction with some of the optimization algorithms to improve convergence and produce interesting future results [1]. Thus, they would be an interesting addition to OCOBox.

50 Conclusion

We introduced a new MATLAB toolbox for online convex optimization. It covers various aspects of the OCO including Boosting, momentum based optimization, and bandit learning. Furthermore, it includes some algorithms on learning rate scheduling and interactive machine learning. This toolbox could hopefully grow into an open sourced community in the future, as we accrue more contributors starting at the University of Virginia. Currently, the toolbox is far from complete, and one of the first major things we plan to release is a GUI for the many algorithms included for a better, user-friendly experience. In the future, we plan to also develop a better sense of modularity between algorithms and schedulers, and add a better visualization support system. The current toolbox can be found at: <https://github.com/PhilipNaveen/OCOBox>

References

- [1] Jagdeep Bhatia. *Simple and Fast Algorithms for Interactive Machine Learning with Random Counter-examples*. 2019. arXiv: 1810.00506 [cs.LG].
- [2] Aaron Defazio et al. *When, Why and How Much? Adaptive Learning Rate Scheduling by Refinement*. 2023. arXiv: 2310.07831 [cs.LG].
- [3] Shaddin Dughmi and Haifeng Xu. *Algorithmic Bayesian Persuasion*. 2016. arXiv: 1503.05988 [cs.GT].
- [4] Elad Hazan. “Introduction to Online Convex Optimization”. In: *Found. Trends Optim.* 2 (2016), pp. 157–325. URL: <https://api.semanticscholar.org/CorpusID:30482768>.
- [5] Klaus Höllig and Jörg Hörner. “Matlab®”. In: *Aufgaben und Lösungen zur Höheren Mathematik 1* (2020). URL: <https://api.semanticscholar.org/CorpusID:239294939>.
- [6] Amri Omar et al. “The commonly used algorithms to optimize a neural network in supervised learning: Overview, and comparative study”. In: *2021 International Conference on Digital Age & Technological Advances for Sustainable Development (ICDATA)* (2021), pp. 31–38. URL: <https://api.semanticscholar.org/CorpusID:243946835>.
- [7] Warren B Powell. “Reinforcement Learning and Stochastic Optimization”. In: *Reinforcement Learning and Stochastic Optimization: A Unified Framework for Sequential Decisions* (2022). URL: <https://api.semanticscholar.org/CorpusID:235632698>.
- [8] Shai Shalev-Shwartz. “Online Learning and Online Convex Optimization”. In: *Found. Trends Mach. Learn.* 4 (2012), pp. 107–194. URL: <https://api.semanticscholar.org/CorpusID:51730029>.
- [9] Li Yang and Abdallah Shami. “On Hyperparameter Optimization of Machine Learning Algorithms: Theory and Practice”. In: *ArXiv abs/2007.15745* (2020). URL: <https://api.semanticscholar.org/CorpusID:220919678>.