

Parameter-free Online Optimization Part 1

Francesco Orabona Ashok Cutkosky

ICML 2020

- Most of the algorithms we know/use in machine learning have parameters

- Most of the algorithms we know/use in machine learning have parameters
 - E.g. regularization parameter in LASSO, k in k -Nearest Neighbourhood, topology of the networks in deep learning

- Most of the algorithms we know/use in machine learning have parameters
 - E.g. regularization parameter in LASSO, k in k -Nearest Neighbourhood, topology of the networks in deep learning
- Most of the time, a large enough validation set can be used to tune the parameters

- Most of the algorithms we know/use in machine learning have parameters
 - E.g. regularization parameter in LASSO, k in k -Nearest Neighbourhood, topology of the networks in deep learning
- Most of the time, a large enough validation set can be used to tune the parameters
- But, at what computational price? Are they really necessary?

- Most of the algorithms we know/use in machine learning have parameters
 - E.g. regularization parameter in LASSO, k in k -Nearest Neighbourhood, topology of the networks in deep learning
- Most of the time, a large enough validation set can be used to tune the parameters
- But, at what computational price? Are they really necessary?
- Are you ignoring some computational/sample complexities?

- Most of the algorithms we know/use in machine learning have parameters
 - E.g. regularization parameter in LASSO, k in k -Nearest Neighbourhood, topology of the networks in deep learning
- Most of the time, a large enough validation set can be used to tune the parameters
- But, at what computational price? Are they really necessary?
- Are you ignoring some computational/sample complexities?

- When is the last time you needed to tune the learning rate to invert a matrix?

How ML Should Be

- No parameters to tune
- No humans in the loop
- Some guarantees

How ML Should Be

- No parameters to tune
 - No humans in the loop
 - Some guarantees
-
- Sometimes this is possible
 - This tutorial is about a method to achieve it in many interesting cases

Learning Rates are Unavoidable?

Learning Rates are Unavoidable?

- Let's consider an example
- You have 11 real numbers y_1, \dots, y_{11}
- We want to solve $\min_x \sum_{i=1}^{11} |x - y_i|$
- You might use SGD
- It will converge, but you have to set a learning rate...

Learning Rates are Unavoidable?

- Let's consider an example
 - You have 11 real numbers y_1, \dots, y_{11}
 - We want to solve $\min_x \sum_{i=1}^{11} |x - y_i|$
 - You might use SGD
 - It will converge, but you have to set a learning rate...
-
- Alternative: you just need the median, so sort the numbers and pick the one in the middle...

Outline of the Tutorial

- Part 1: Stochastic and Online Convex Optimization
- Part 2: Parameter-free Convex Optimization
- Part 3: More Adaptivity and Applications
- Part 4: Implementation, Experiments, Open Problems

Learning Objectives

- Explain the problems behind learning rates
- Intuition behind recent parameter-free approaches
- Stripped-down algorithms and proofs
- Trick&tips for practical implementations

Overall, democratize these new methods and expand the number of people working in this field

- **Part 1: Stochastic and Online Convex Optimization**
- Part 2: Parameter-free Convex Optimization
- Part 3: More Adaptivity and Applications
- Part 4: Implementation, Experiments, Open Problems

- We want to solve the problem

$$\min_{\mathbf{x} \in X} F(\mathbf{x}),$$

where F is a convex function

- We will also assume F to be G -Lipschitz, i.e., $|F(\mathbf{x}) - F(\mathbf{y})| \leq G\|\mathbf{x} - \mathbf{y}\|$
- For the most part of this tutorial, we will not assume much more than this

- We have access to an first-order stochastic oracle
 - It gives us noisy estimate \mathbf{g}_t of the gradient of F in a point \mathbf{x}_t
 - We will assume unbiasedness: $\mathbb{E}[\mathbf{g}_t | \mathbf{g}_1, \dots, \mathbf{g}_{t-1}] = \nabla F(\mathbf{x}_t)$
- ⚠ Most of the things we say work also for non-differentiable functions using subgradients

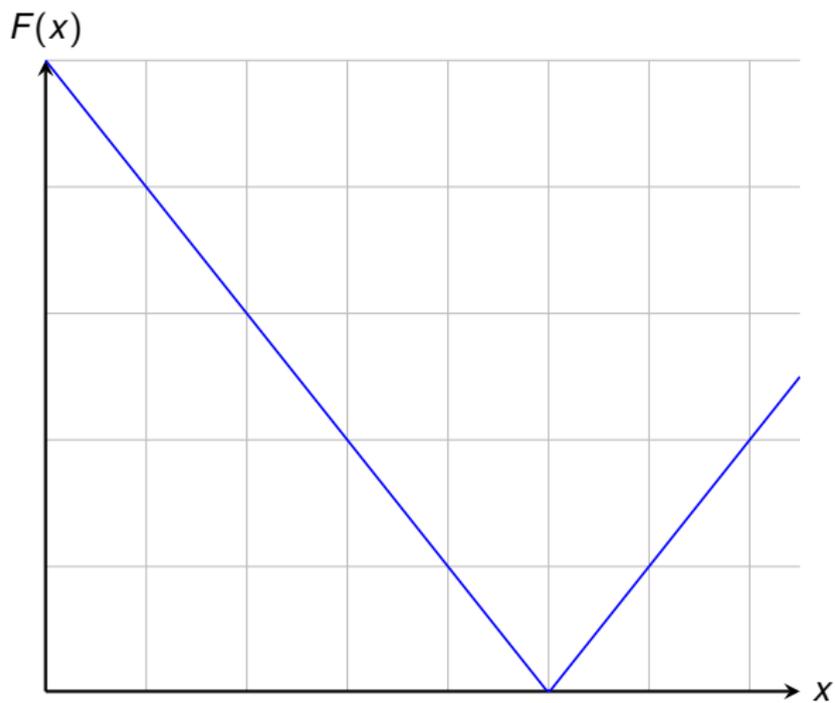
The Most Used Algorithm in ML: SGD

- The simplest algorithm that we can use in our setting is Stochastic Gradient Descent (SGD)
- Iteratively it moves in the negative direction of the gradient:

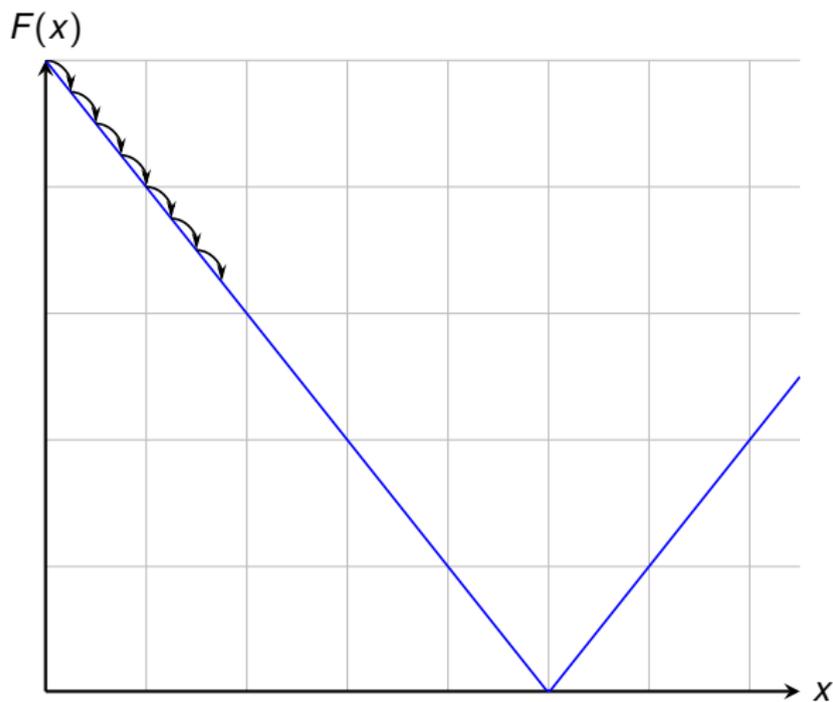
$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \mathbf{g}_t$$

- η_t is the learning rate or step size
- Choosing the correct learning rate is a black art...
- Why is it so difficult?

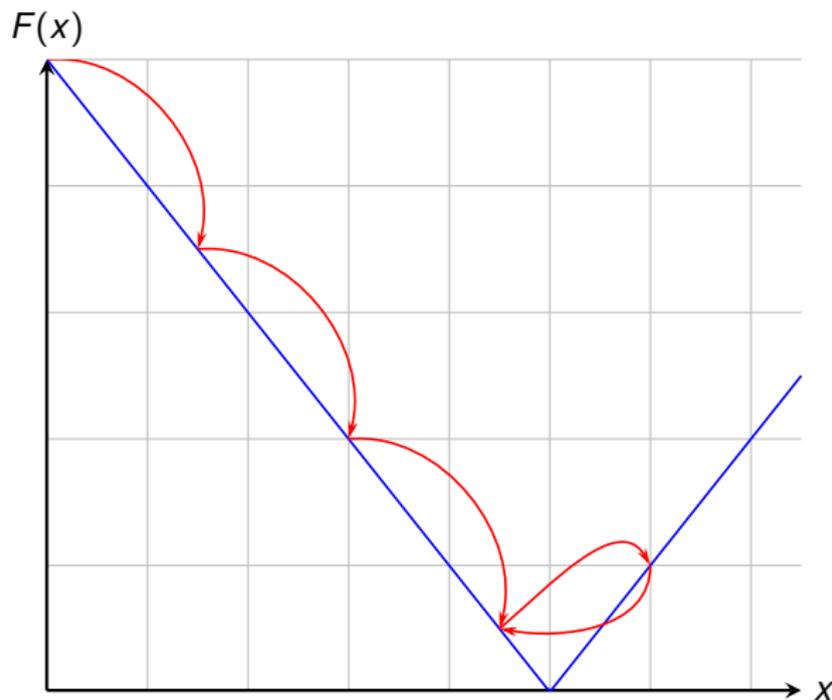
Minimize $F(x) = |x - 10|$



First Problem: The Learning Rate Sets Your Maximum Precision

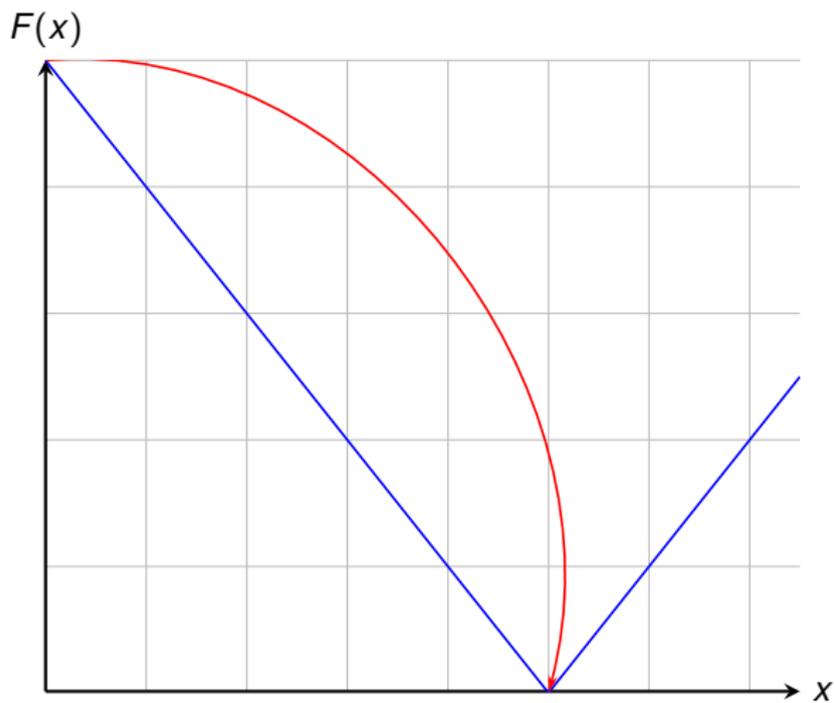


First Problem: The Learning Rate Sets Your Maximum Precision



⚠ If you want to converge at least in the limit, learning rates cannot be fixed

One Jump?

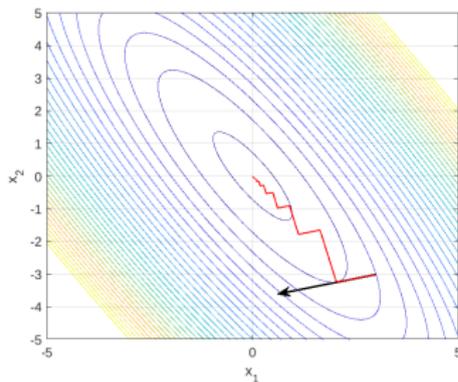
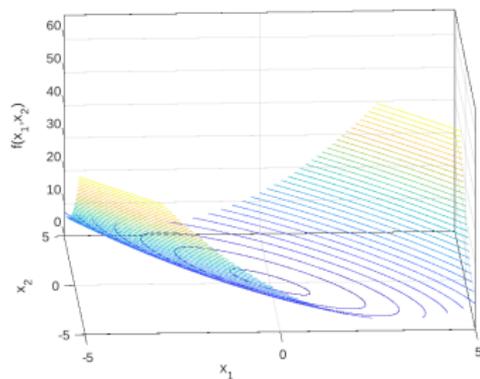


No One Jump!

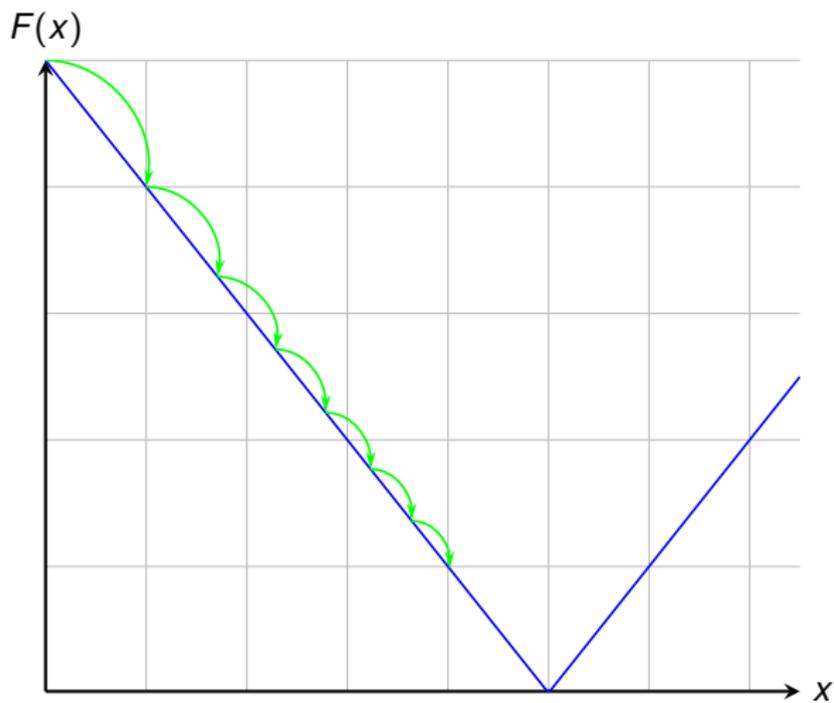
- First, in the stochastic setting, line search is not possible because you do not have the exact gradient nor can you exactly compute the function value
 - You need strong assumptions (and parameters...) to make stochastic line search work [Paquette&Scheinberg, arXiv'18; Vaswani et al., NeurIPS'19]

No One Jump!

- First, in the stochastic setting, line search is not possible because you do not have the exact gradient nor can you exactly compute the function value
 - You need strong assumptions (and parameters...) to make stochastic line search work [Paquette&Scheinberg, arXiv'18; Vaswani et al., NeurIPS'19]
- Even in 2d, the gradients does not necessarily point towards the minimum

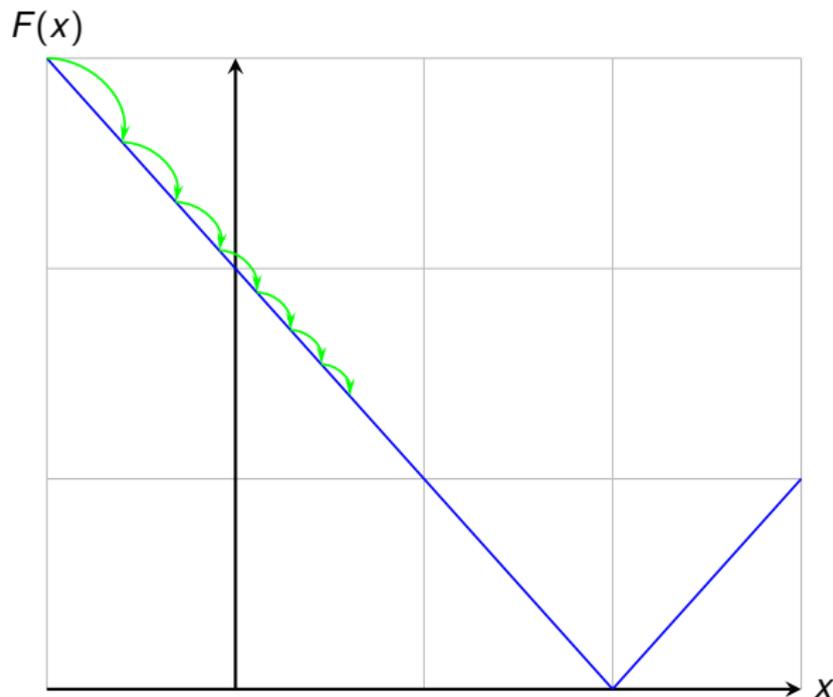


Solution: Decreasing Learning Rates



Second Problem: Distance to the Optimum

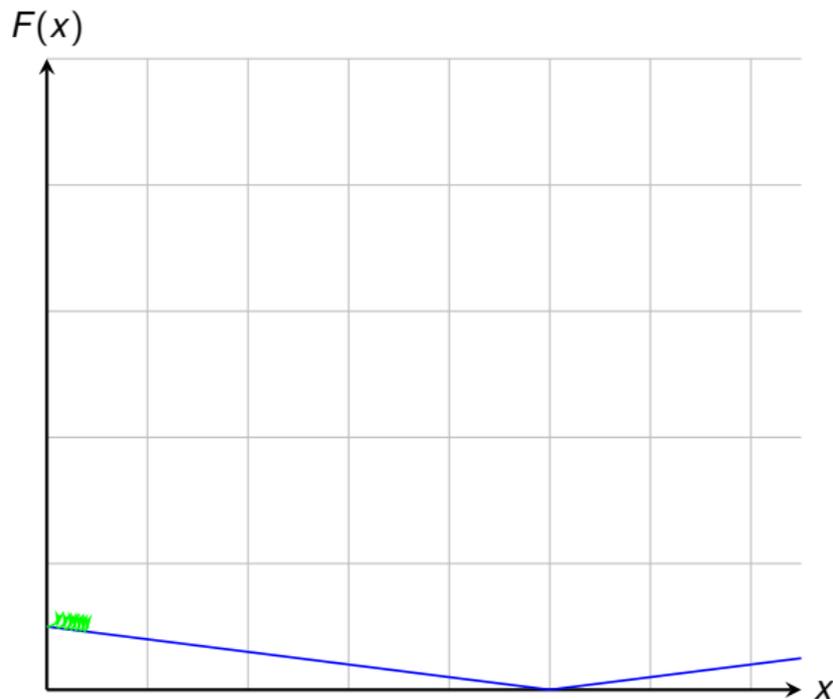
Starting from -5, our choice of the learning rate is not good anymore!



⚠️ Optimal learning rate depends on distance from the optimum

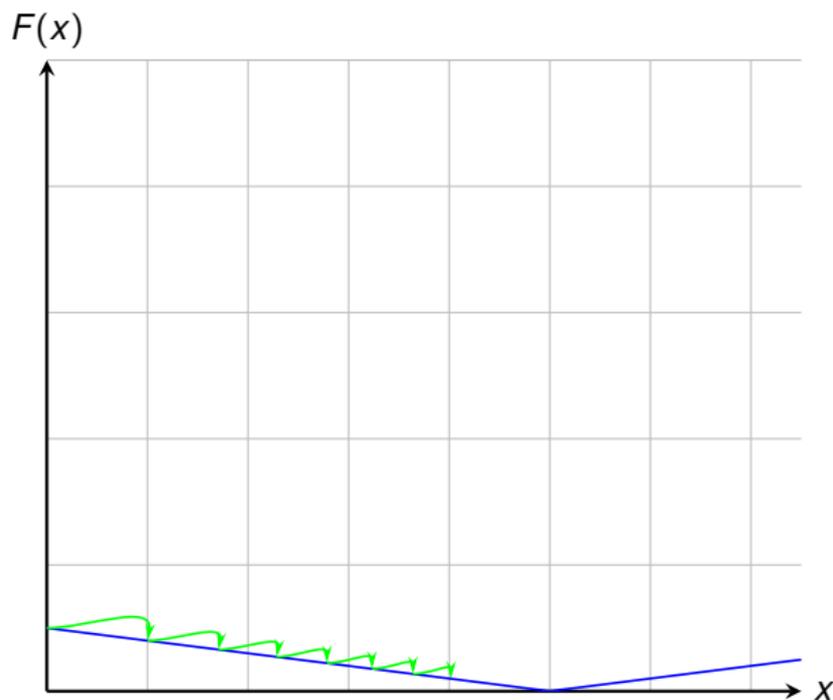
Third Problem: Scale

Changing the function to $0.1|x - 10|$ again ruins the convergence



⚠ Optimal learning rate depends on the scale (=Lipschitz constant)

What About “Adaptive” Algorithms?



- ⚠ AdaGrad adapts to the scale of each coordinate, but does not adapt to distance from the optimum!

Take Home Messages on Learning Rates

Learning rates...

- 1 ...must decrease
- 2 ...depend on the distance to the optimal solution
- 3 ...depend on the scale

The Theory Confirms our Intuition

- Use a decreasing learning rate, $\eta_t = \frac{\alpha}{\sqrt{t}}$
- Convergence after T iterations is $\mathcal{O}\left(\frac{1}{\sqrt{T}}\left(\frac{\|\mathbf{x}^* - \mathbf{x}_1\|^2}{\alpha} + \alpha G^2\right)\right)$
- \mathbf{x}^* is the best solution
- G the bound on the norm of the stochastic gradients

The Theory Confirms our Intuition

- Use a decreasing learning rate, $\eta_t = \frac{\alpha}{\sqrt{t}}$
- Convergence after T iterations is $\mathcal{O}\left(\frac{1}{\sqrt{T}}\left(\frac{\|\mathbf{x}^* - \mathbf{x}_1\|^2}{\alpha} + \alpha G^2\right)\right)$
- \mathbf{x}^* is the best solution
- G the bound on the norm of the stochastic gradients

- The optimal learning rate is with $\alpha = \frac{\|\mathbf{x}^* - \mathbf{x}_1\|}{G}$ that would give a rate of $\mathcal{O}\left(\frac{G\|\mathbf{x}^* - \mathbf{x}_1\|}{\sqrt{T}}\right)$

The Theory Confirms our Intuition

- Use a decreasing learning rate, $\eta_t = \frac{\alpha}{\sqrt{t}}$
- Convergence after T iterations is $\mathcal{O}\left(\frac{1}{\sqrt{T}}\left(\frac{\|\mathbf{x}^* - \mathbf{x}_1\|^2}{\alpha} + \alpha G^2\right)\right)$
- \mathbf{x}^* is the best solution
- G the bound on the norm of the stochastic gradients

- The optimal learning rate is with $\alpha = \frac{\|\mathbf{x}^* - \mathbf{x}_1\|}{G}$ that would give a rate of $\mathcal{O}\left(\frac{G\|\mathbf{x}^* - \mathbf{x}_1\|}{\sqrt{T}}\right)$
- ...but you don't know \mathbf{x}^* ...

The Theory Confirms our Intuition

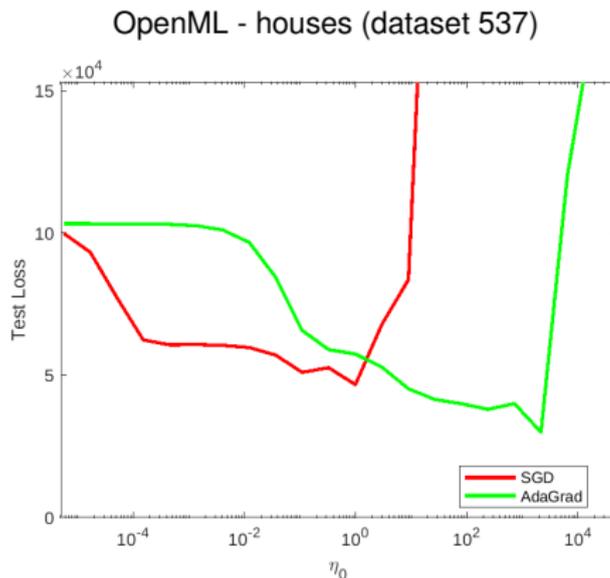
- Use a decreasing learning rate, $\eta_t = \frac{\alpha}{\sqrt{t}}$
- Convergence after T iterations is $\mathcal{O}\left(\frac{1}{\sqrt{T}}\left(\frac{\|\mathbf{x}^* - \mathbf{x}_1\|^2}{\alpha} + \alpha G^2\right)\right)$
- \mathbf{x}^* is the best solution
- G the bound on the norm of the stochastic gradients

- The optimal learning rate is with $\alpha = \frac{\|\mathbf{x}^* - \mathbf{x}_1\|}{G}$ that would give a rate of $\mathcal{O}\left(\frac{G\|\mathbf{x}^* - \mathbf{x}_1\|}{\sqrt{T}}\right)$
- ...but you don't know \mathbf{x}^* ...

- Only if the domain is bounded, you can (loosely) upper bound $\|\mathbf{x}^* - \mathbf{x}_1\|$
 - Not-a-solution solution: Assume your domain is bounded or bounded iterates

I Hate Learning Rates

- Finding the right one is critical!
- The only existing solution: tries different learning rates



Theorem

- $\mathcal{O}\left(\frac{G\|\mathbf{x}^* - \mathbf{x}_1\|}{\sqrt{T}}\right)$ *cannot be achieved, unless we know $\|\mathbf{x}^* - \mathbf{x}_1\|$*

Theorem

- $\mathcal{O}\left(\frac{G\|\mathbf{x}^* - \mathbf{x}_1\|}{\sqrt{T}}\right)$ *cannot* be achieved, unless we know $\|\mathbf{x}^* - \mathbf{x}_1\|$
- we can achieve $\mathcal{O}\left(\frac{G\|\mathbf{x}^* - \mathbf{x}_1\|}{\sqrt{T}} \sqrt{\log(G\|\mathbf{x}^* - \mathbf{x}_1\|\sqrt{T} + 1)} + \frac{1}{T}\right)$

Theorem

- $\mathcal{O}\left(\frac{G\|\mathbf{x}^* - \mathbf{x}_1\|}{\sqrt{T}}\right)$ *cannot* be achieved, unless we know $\|\mathbf{x}^* - \mathbf{x}_1\|$
- we can achieve $\mathcal{O}\left(\frac{G\|\mathbf{x}^* - \mathbf{x}_1\|}{\sqrt{T}} \sqrt{\log(G\|\mathbf{x}^* - \mathbf{x}_1\|\sqrt{T} + 1)} + \frac{1}{T}\right)$
- or $\mathcal{O}\left(\frac{G\|\mathbf{x}^* - \mathbf{x}_1\|}{\sqrt{T}} \sqrt{\log(G\|\mathbf{x}^* - \mathbf{x}_1\| + 1)} + \frac{1}{\sqrt{T}}\right)$

Theorem

- $\mathcal{O}\left(\frac{G\|\mathbf{x}^* - \mathbf{x}_1\|}{\sqrt{T}}\right)$ *cannot* be achieved, unless we know $\|\mathbf{x}^* - \mathbf{x}_1\|$
 - we can achieve $\mathcal{O}\left(\frac{G\|\mathbf{x}^* - \mathbf{x}_1\|}{\sqrt{T}} \sqrt{\log(G\|\mathbf{x}^* - \mathbf{x}_1\|\sqrt{T} + 1)} + \frac{1}{T}\right)$
 - or $\mathcal{O}\left(\frac{G\|\mathbf{x}^* - \mathbf{x}_1\|}{\sqrt{T}} \sqrt{\log(G\|\mathbf{x}^* - \mathbf{x}_1\| + 1)} + \frac{1}{\sqrt{T}}\right)$
-
- The logarithmic term is the price you have to pay because you don't know how far is the optimal solution
 - Think of it as the complexity of binary searching for the optimal learning rate

- Before solving this problem, we need some tools
- We will use Online Convex Optimization algorithms

- Online Optimization \neq Stochastic Optimization
- Online Optimization is actually more general than Stochastic Optimization!

Online Learning

- 1 In each round, output $\mathbf{x}_t \in X$
- 2 Pay $\ell_t(\mathbf{x}_t)$
- 3 Receive $\mathbf{g}_t = \nabla \ell_t(\mathbf{x}_t)$
- 4 Update \mathbf{x}_{t+1}

From Online Learning to Stochastic Optimization

Online Learning

- 1 In each round, output $\mathbf{x}_t \in \mathcal{X}$
- 2 Pay $\ell_t(\mathbf{x}_t)$
- 3 Receive $\mathbf{g}_t = \nabla \ell_t(\mathbf{x}_t)$
- 4 Update \mathbf{x}_{t+1}

Choose \mathbf{x}_t before observing ℓ_t

From Online Learning to Stochastic Optimization

Online Learning

- 1 In each round, output $\mathbf{x}_t \in X$
- 2 Pay $\ell_t(\mathbf{x}_t)$
- 3 Receive $\mathbf{g}_t = \nabla \ell_t(\mathbf{x}_t)$
- 4 Update \mathbf{x}_{t+1}

Choose \mathbf{x}_t before observing ℓ_t
No assumptions on how ℓ_t is generated!

From Online Learning to Stochastic Optimization

Online Learning

- 1 In each round, output $\mathbf{x}_t \in X$
- 2 Pay $\ell_t(\mathbf{x}_t)$
- 3 Receive $\mathbf{g}_t = \nabla \ell_t(\mathbf{x}_t)$
- 4 Update \mathbf{x}_{t+1}

Choose \mathbf{x}_t before observing ℓ_t
No assumptions on how ℓ_t is generated!

Regret minimization

$$\min_{\mathbf{x}_1, \dots, \mathbf{x}_T \in X} \sum_{t=1}^T \ell_t(\mathbf{x}_t) \quad \text{equivalently} \quad \min_{\mathbf{x}_1, \dots, \mathbf{x}_T \in X} \underbrace{\sum_{t=1}^T \ell_t(\mathbf{x}_t) - \sum_{t=1}^T \ell_t(\mathbf{x}^*)}_{\text{Regret}_T(\mathbf{x}^*)}$$

- Typically, we can design Online algorithms where the $\text{Regret}_T(\mathbf{x}^*) = O(\sqrt{T})$
- Harder setting, easier analysis, same results :)

From Online to Stochastic (or Batch) Optimization

- 1: **for** $t = 1$ **to** T **do**
- 2: Get \mathbf{x}_t from Online Learning Algorithm
- 3: Receive stochastic gradient \mathbf{g}_t such that $\mathbb{E}_t[\mathbf{g}_t] = \nabla F(\mathbf{x}_t)$
- 4: Pass loss $\ell_t(\mathbf{x}) = \langle \mathbf{g}_t, \mathbf{x} \rangle$ to Online Learning Algorithm
- 5: **end for**
- 6: **return** $\bar{\mathbf{x}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t$

Theorem

$$\mathbb{E}[F(\bar{\mathbf{x}}_T)] - F(\mathbf{x}^*) \leq \frac{\mathbb{E}[\text{Regret}_T(\mathbf{x}^*)]}{T}$$

In words: If you have an online algorithm able to work just with linear losses, you have a stochastic optimization algorithm for any convex function!

Some Famous Online Learning Algorithms

- Online Gradient Descent [Zinkevich, ICML'03]
- AdaGrad [Duchi et al., COLT'10, JMRL'11]
- AMSGrad [Reddi et al., ICLR'18]

These algorithms are designed to work in the adversarial setting and have a $O(\sqrt{T})$ regret bound

Hence, they can also be used as stochastic optimization algorithms with a $O(\frac{1}{\sqrt{T}})$ convergence rate

Summary of Part 1

- 1 Learning rates are difficult to set because they depends on quantities you don't know nor you can easily estimate
- 2 In particular, the additional logarithmic dependency on the distance to the optimal solution is unavoidable
- 3 Online Optimization gives us a set of powerful tools to attack this problem