

Parameter-free Online Optimization Part 2

Francesco Orabona Ashok Cutkosky

ICML 2020

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

- We saw why tuning learning rates is difficult
- In particular, the distance from to the optimum is unknown and it cannot be easily estimated
- We want to achieve $\mathcal{O}\left(\frac{G\|\mathbf{x}^* - \mathbf{x}_1\|}{\sqrt{T}} \sqrt{\log(G\|\mathbf{x}^* - \mathbf{x}_1\|\sqrt{T} + 1)} + \frac{1}{T}\right)$

- We saw why tuning learning rates is difficult
- In particular, the distance from to the optimum is unknown and it cannot be easily estimated
- We want to achieve $\mathcal{O}\left(\frac{G\|\mathbf{x}^* - \mathbf{x}_1\|}{\sqrt{T}} \sqrt{\log(G\|\mathbf{x}^* - \mathbf{x}_1\|\sqrt{T} + 1)} + \frac{1}{T}\right)$
- What next?

How Would Nemirovski Remove the Learning Rate? (1)

Idea: Try a grid of learning rates, but with a budget for each one, so that overall you spend the same time as using only 1 learning rate.

- Procedure $M(\alpha, T)$: Run GD for T iterations, with stepsizes $\eta_t = \frac{\alpha}{\sqrt{T} \|\nabla F(\mathbf{x}_t)\|}$, starting at the origin, and return $\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t$
- For $n = 1, 2, \dots$ let us run one after one the procedures $M(\alpha = 1, T = 2^n), M(\alpha = 2, T = 2^n), M(\alpha = 4, T = 2^n), M(\alpha = 8, T = 2^n), \dots, M(\alpha = 2^n, T = 2^n)$ and return the one with the smallest value of F
- Recall that we need $\alpha \approx \|\mathbf{x}^*\|$
- Set m be such that $2^{m-1} \leq \|\mathbf{x}^*\| \leq 2^m$

How Would Nemirovski Remove the Learning Rate? (2)

- When $n \geq m$, the n -th batch will produce a solution of accuracy $O(1)G\|\mathbf{x}^*\|2^{-n/2}$ – this is what you will get when running $M(\alpha = 2^m, T = 2^n)$
- How long does it take?
- The total number of steps, from the beginning till executing n -th batch will be $T(n) = O(1)n2^n$
- So, when the total number of steps T satisfies $T \geq O(1)m2^m = O(1)\|\mathbf{x}^*\| \ln(\|\mathbf{x}^*\| + 1)$, the inaccuracy will be at most $O(1)\frac{\|\mathbf{x}^*\|G\sqrt{\ln T}}{\sqrt{T}}$

Not a Satisfying Solution...

- ✓ The construction tells us in principle there is an algorithm that finds the optimal learning rate
- × It does not work in the stochastic setting: one of the runs was good, but you don't know which one
- × You “waste” computation trying wrong learning rates
- × It does not work in practice

Second Alternative: Doubling Trick

Require: Initial Learning rate η_0

```
1:  $x_1 = 0, Q_0 = 0$ 
2: for  $t = 1, \dots, T$  do
3:   Output  $x_t$ 
4:   Receive gradient  $g_t = \nabla F(x_t)$ 
5:    $Q_t = Q_{t-1} - g_t x_t$ 
6:   if  $Q_t \leq \eta_{t-1} T$  then
7:      $\eta_t = \eta_{t-1}$ 
8:   else
9:      $x_t = 0$ 
10:     $\eta_t = 2\eta_{t-1}$ 
11:     $Q_t = 0$ 
12:   end if
13:    $x_{t+1} = x_t - \eta_t g_t$ 
14: end for
```

- ✓ It works in the stochastic and online setting
- ✗ The bound is suboptimal
- ✗ It does not really work (due to the resets)

Third Alternative: Multiple Learning Rates in “Parallel”

- Instead of trying each single learning rate on its own, we might try all of them at the same time
 - We run $N = O(\log T)$ SGD procedures, each with a different learning rate
 - We feed the N algorithms with surrogate gradients and we merge their outputs at each step
 - This approach was first proposed in the MetaGrad [van Erven&Koolen, NeurIPS'16] and Maler [Wang et al., UAI'19]
 - Foster et al. [NeurIPS'17] proposed a similar approach that works for unbounded domains too
-
- ✓ General approach: you can use any base optimizer, not only SGD
 - × MetaGrad and Maler work only for bounded domains
 - × These approaches are a bit cumbersome: you do need to run all the algorithms at the same time

Fourth Alternative: FTRL/RDA Approach

- We would like to obtain the optimal convergence rate

$$O\left(\frac{\|\mathbf{x}^* - \mathbf{x}_1\| \sqrt{\ln(\|\mathbf{x}^* - \mathbf{x}_1\| \sqrt{T})}}{\sqrt{T}} + \frac{1}{T}\right)$$

- Equivalently, we would like to obtain the optimal regret

$$\text{Regret}_T(\mathbf{u}) \leq O\left(\|\mathbf{u}\| \sqrt{T \ln(\|\mathbf{u}\| \sqrt{T})} + 1\right)$$

- So, we can use Follow-the-Regularized-Leader (aka Regularized Dual Averaging) with time-varying regularizer $\psi_t(\mathbf{x}) \approx \|\mathbf{x}\| \sqrt{t \ln(\|\mathbf{x}\| \sqrt{t})}$

$$\mathbf{x}_t \approx \underset{\mathbf{x}}{\text{argmin}} \|\mathbf{x}\| \sqrt{t \ln(\|\mathbf{x}\| \sqrt{t})} + \sum_{i=1}^{t-1} \langle \mathbf{x}, \mathbf{g}_i \rangle$$

- ✓ It works!
- × Very difficult to analyze and explain

Last Alternative: Coin-betting Approach

This is the same as the FTRL approach, but:

- ✓ Easy
- ✓ Gives a new interpretation to optimization as gambling/compression/prediction with log loss
- ✓ Works very well :)
- × Very different from the usual optimization algorithms

Betting on a Coin



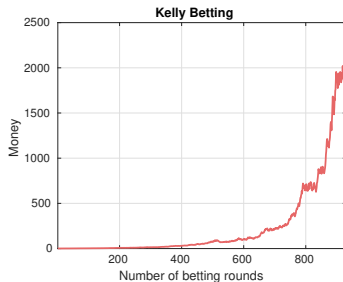
- Start with \$1
- Bet x_t money on head ($x_t > 0$) or tails ($x_t < 0$)
 - Cannot borrow money
- Win or lose depending on the outcome of the coin
 $c_t \in \{-1, 1\}$
- $Wealth_t = Wealth_{t-1} + x_t c_t$



Aim: Maximize gain on all sequences where the number of tails and head are different

Optimal Betting Strategy for a Stochastic Coin: Kelly Betting (1956)

- Known problem in economics.
- Bet a **fraction** of your **money** equal to $2p - 1 = \mathbb{E}[c_t]$ on tail at each round.
 - E.g. $p = 0.51 \Rightarrow$ bet **2%**
- Expected log wealth is linear in time.



Non-Stochastic setting, but Knowing the Future

- Non-stochastic setting, T rounds.
- Assume we bet a fixed fraction of money at each round.
- What is the optimal fraction?

Non-Stochastic setting, but Knowing the Future

- Non-stochastic setting, T rounds.
- Assume we bet a fixed fraction of money at each round.
- What is the optimal fraction?
- The optimal signed fraction at each time step t is $\frac{\sum_{i=1}^T c_i}{T}$.
- Hence you bet $\frac{\sum_{i=1}^T c_i}{T} \cdot \text{Wealth}_{t-1}$.
- Winnings are exponential:

$$\text{Winnings} > \exp\left(\frac{(\sum_{t=1}^T c_t)^2}{2T}\right)$$

Worst case Optimal Betting Strategy: Krichevsky-Trofimov Better

- What if we don't know the future?

Worst case Optimal Betting Strategy: Krichevsky-Trofimov Bettor

- What if we don't know the future?
- Estimate “probability” of heads with Krichevsky-Trofimov (KT) estimator:

$$\frac{\frac{1}{2} + \# \text{ of heads in } t \text{ rounds}}{t+1}$$

- No stochastic assumptions: KT has optimal regret w.r.t. the log loss
- Not the only possible strategy, but the simplest one

Worst case Optimal Betting Strategy: Krichevsky-Trofimov Bettor

- What if we don't know the future?
- Estimate “probability” of heads with Krichevsky-Trofimov (KT) estimator:
$$\frac{\frac{1}{2} + \# \text{ of heads in } t \text{ rounds}}{t+1}$$
 - No stochastic assumptions: KT has optimal regret w.r.t. the log loss
 - Not the only possible strategy, but the simplest one
- Hence, on round t bet a signed fraction of your money equal to $\frac{\sum_{i=1}^{t-1} c_i}{t}$

Worst case Optimal Betting Strategy: Krichevsky-Trofimov Bettor

- What if we don't know the future?
- Estimate “probability” of heads with Krichevsky-Trofimov (KT) estimator:
$$\frac{\frac{1}{2} + \# \text{ of heads in } t \text{ rounds}}{t+1}$$
 - No stochastic assumptions: KT has optimal regret w.r.t. the log loss
 - Not the only possible strategy, but the simplest one
- Hence, on round t bet a signed fraction of your money equal to $\frac{\sum_{i=1}^{t-1} c_i}{t}$
- Still exponential

$$\text{Winnings of KT Bettor} \geq \frac{\exp\left(\frac{(\sum_{t=1}^T c_t)^2}{2T}\right)}{2\sqrt{T}}$$

Worst case Optimal Betting Strategy: Krichevsky-Trofimov Bettor

- What if we don't know the future?
- Estimate “probability” of heads with Krichevsky-Trofimov (KT) estimator:
$$\frac{\frac{1}{2} + \# \text{ of heads in } t \text{ rounds}}{t+1}$$
 - No stochastic assumptions: KT has optimal regret w.r.t. the log loss
 - Not the only possible strategy, but the simplest one
- Hence, on round t bet a signed fraction of your money equal to $\frac{\sum_{i=1}^{t-1} c_i}{t}$
- Still exponential

$$\text{Winnings of KT Bettor} \geq \frac{\exp\left(\frac{(\sum_{t=1}^T c_t)^2}{2T}\right)}{2\sqrt{T}} = \frac{\text{Winnings knowing the future}}{2\sqrt{T}}$$

Worst case Optimal Betting Strategy: Krichevsky-Trofimov Bettor

- What if we don't know the future?
- Estimate “probability” of heads with Krichevsky-Trofimov (KT) estimator:
$$\frac{1}{2} + \frac{\text{\# of heads in } t \text{ rounds}}{t+1}$$
 - No stochastic assumptions: KT has optimal regret w.r.t. the log loss
 - Not the only possible strategy, but the simplest one
- Hence, on round t bet a signed fraction of your money equal to $\frac{\sum_{i=1}^{t-1} c_i}{t}$
- Still exponential

$$\text{Winnings of KT Bettor} \geq \frac{\exp\left(\frac{(\sum_{t=1}^T c_t)^2}{2T}\right)}{2\sqrt{T}} = \frac{\text{Winnings knowing the future}}{2\sqrt{T}}$$

Coin-betting is solvable with a parameter-free optimal strategy, but what is the connection with SGD?

1D Optimization through Betting

- Assume that there exists a function $H(\cdot)$ such that our betting strategy will guarantee that the wealth after T rounds will be at least $H(\sum_{t=1}^T c_t)$ for any arbitrary sequence c_1, \dots, c_T

1D Optimization through Betting

- Assume that there exists a function $H(\cdot)$ such that our betting strategy will guarantee that the wealth after T rounds will be at least $H(\sum_{t=1}^T c_t)$ for any arbitrary sequence c_1, \dots, c_T
 - For KT $H(\sum_{t=1}^T c_t) = C \exp((\sum_{t=1}^T c_t)^2 / T) / \sqrt{T}$

1D Optimization through Betting

- Assume that there exists a function $H(\cdot)$ such that our betting strategy will guarantee that the wealth after T rounds will be at least $H(\sum_{t=1}^T c_t)$ for any arbitrary sequence c_1, \dots, c_T
 - For KT $H(\sum_{t=1}^T c_t) = C \exp((\sum_{t=1}^T c_t)^2 / T) / \sqrt{T}$
- We want to minimize $F(x) = |x - 10|$

1D Optimization through Betting

- Assume that there exists a function $H(\cdot)$ such that our betting strategy will guarantee that the wealth after T rounds will be at least $H(\sum_{t=1}^T c_t)$ for any arbitrary sequence c_1, \dots, c_T
 - For KT $H(\sum_{t=1}^T c_t) = C \exp((\sum_{t=1}^T c_t)^2 / T) / \sqrt{T}$
- We want to minimize $F(x) = |x - 10|$
- Let's set a betting game: bet x_t dollars on $c_t = -g_t = -\nabla F(x_t) \in \{-1, +1\}$

1D Optimization through Betting

- Assume that there exists a function $H(\cdot)$ such that our betting strategy will guarantee that the wealth after T rounds will be at least $H(\sum_{t=1}^T c_t)$ for any arbitrary sequence c_1, \dots, c_T
 - For KT $H(\sum_{t=1}^T c_t) = C \exp((\sum_{t=1}^T c_t)^2 / T) / \sqrt{T}$
- We want to minimize $F(x) = |x - 10|$
- Let's set a betting game: bet x_t dollars on $c_t = -g_t = -\nabla F(x_t) \in \{-1, +1\}$
- Claim: the average of the bets will converge to the minimum of $F(x)$ at a rate that depends on how good is our betting strategy!

$$F\left(\frac{1}{T} \sum_{t=1}^T x_t\right) - F(x^*)$$

1D Optimization through Betting

- Assume that there exists a function $H(\cdot)$ such that our betting strategy will guarantee that the wealth after T rounds will be at least $H(\sum_{t=1}^T c_t)$ for any arbitrary sequence c_1, \dots, c_T
 - For KT $H(\sum_{t=1}^T c_t) = C \exp((\sum_{t=1}^T c_t)^2 / T) / \sqrt{T}$
- We want to minimize $F(x) = |x - 10|$
- Let's set a betting game: bet x_t dollars on $c_t = -g_t = -\nabla F(x_t) \in \{-1, +1\}$
- Claim: the average of the bets will converge to the minimum of $F(x)$ at a rate that depends on how good is our betting strategy!

$$F\left(\frac{1}{T} \sum_{t=1}^T x_t\right) - F(x^*) \stackrel{\text{Jensen}}{\leq} \frac{1}{T} \sum_{t=1}^T F(x_t) - F(x^*)$$

1D Optimization through Betting

- Assume that there exists a function $H(\cdot)$ such that our betting strategy will guarantee that the wealth after T rounds will be at least $H(\sum_{t=1}^T c_t)$ for any arbitrary sequence c_1, \dots, c_T
 - For KT $H(\sum_{t=1}^T c_t) = C \exp((\sum_{t=1}^T c_t)^2 / T) / \sqrt{T}$
- We want to minimize $F(x) = |x - 10|$
- Let's set a betting game: bet x_t dollars on $c_t = -g_t = -\nabla F(x_t) \in \{-1, +1\}$
- Claim: the average of the bets will converge to the minimum of $F(x)$ at a rate that depends on how good is our betting strategy!

$$F\left(\frac{1}{T} \sum_{t=1}^T x_t\right) - F(x^*) \stackrel{\text{Jensen}}{\leq} \frac{1}{T} \sum_{t=1}^T F(x_t) - F(x^*) \stackrel{\text{Convexity}}{\leq} \frac{1}{T} \left(\sum_{t=1}^T c_t x^* - \sum_{t=1}^T c_t x_t \right)$$

1D Optimization through Betting

- Assume that there exists a function $H(\cdot)$ such that our betting strategy will guarantee that the wealth after T rounds will be at least $H(\sum_{t=1}^T c_t)$ for any arbitrary sequence c_1, \dots, c_T
 - For KT $H(\sum_{t=1}^T c_t) = C \exp((\sum_{t=1}^T c_t)^2 / T) / \sqrt{T}$
- We want to minimize $F(x) = |x - 10|$
- Let's set a betting game: bet x_t dollars on $c_t = -g_t = -\nabla F(x_t) \in \{-1, +1\}$
- Claim: the average of the bets will converge to the minimum of $F(x)$ at a rate that depends on how good is our betting strategy!

$$\begin{aligned} F\left(\frac{1}{T} \sum_{t=1}^T x_t\right) - F(x^*) &\stackrel{\text{Jensen}}{\leq} \frac{1}{T} \sum_{t=1}^T F(x_t) - F(x^*) \stackrel{\text{Convexity}}{\leq} \frac{1}{T} \left(\sum_{t=1}^T c_t x^* - \sum_{t=1}^T c_t x_t \right) \\ &\stackrel{\text{Def. } H}{\leq} \frac{1}{T} \left(\sum_{t=1}^T c_t x^* - H\left(\sum_{t=1}^T c_t\right) + 1 \right) \end{aligned}$$

1D Optimization through Betting

- Assume that there exists a function $H(\cdot)$ such that our betting strategy will guarantee that the wealth after T rounds will be at least $H(\sum_{t=1}^T c_t)$ for any arbitrary sequence c_1, \dots, c_T
 - For KT $H(\sum_{t=1}^T c_t) = C \exp((\sum_{t=1}^T c_t)^2 / T) / \sqrt{T}$
- We want to minimize $F(x) = |x - 10|$
- Let's set a betting game: bet x_t dollars on $c_t = -g_t = -\nabla F(x_t) \in \{-1, +1\}$
- Claim: the average of the bets will converge to the minimum of $F(x)$ at a rate that depends on how good is our betting strategy!

$$\begin{aligned} F\left(\frac{1}{T} \sum_{t=1}^T x_t\right) - F(x^*) &\stackrel{\text{Jensen}}{\leq} \frac{1}{T} \sum_{t=1}^T F(x_t) - F(x^*) \stackrel{\text{Convexity}}{\leq} \frac{1}{T} \left(\sum_{t=1}^T c_t x^* - \sum_{t=1}^T c_t x_t \right) \\ &\stackrel{\text{Def. } H}{\leq} \frac{1}{T} \left(\sum_{t=1}^T c_t x^* - H\left(\sum_{t=1}^T c_t\right) + 1 \right) \stackrel{\text{Max}}{\leq} \frac{1}{T} \left(1 + \max_v v x^* - H(v) \right) \end{aligned}$$

1D Optimization through Betting

- Assume that there exists a function $H(\cdot)$ such that our betting strategy will guarantee that the wealth after T rounds will be at least $H(\sum_{t=1}^T c_t)$ for any arbitrary sequence c_1, \dots, c_T
 - For KT $H(\sum_{t=1}^T c_t) = C \exp((\sum_{t=1}^T c_t)^2 / T) / \sqrt{T}$
- We want to minimize $F(x) = |x - 10|$
- Let's set a betting game: bet x_t dollars on $c_t = -g_t = -\nabla F(x_t) \in \{-1, +1\}$
- Claim: the average of the bets will converge to the minimum of $F(x)$ at a rate that depends on how good is our betting strategy!

$$\begin{aligned} F\left(\frac{1}{T} \sum_{t=1}^T x_t\right) - F(x^*) &\stackrel{\text{Jensen}}{\leq} \frac{1}{T} \sum_{t=1}^T F(x_t) - F(x^*) \stackrel{\text{Convexity}}{\leq} \frac{1}{T} \left(\sum_{t=1}^T c_t x^* - \sum_{t=1}^T c_t x_t \right) \\ &\stackrel{\text{Def. } H}{\leq} \frac{1}{T} \left(\sum_{t=1}^T c_t x^* - H\left(\sum_{t=1}^T c_t\right) + 1 \right) \stackrel{\text{Max}}{\leq} \frac{1}{T} \left(1 + \max_v v x^* - H(v) \right) \stackrel{\text{Def. } H^*}{=} \frac{H^*(x^*) + 1}{T} \end{aligned}$$

KT as an Optimization Algorithm

- Assume $\nabla F(x_t) \in \{-1, +1\}$
- $x_t = \frac{\sum_{i=1}^{t-1} c_i}{t} \text{Wealth}_{t-1} = \frac{-\sum_{i=1}^{t-1} g_i}{t} (1 - \sum_{i=1}^{t-1} g_i \cdot x_i)$

Theorem (Informal)

KT betting in 1-d guarantees

$$F\left(\frac{1}{T} \sum_{t=1}^T x_t\right) - F(x^*) \leq C \frac{|x^*| \sqrt{\log(T)}}{\sqrt{T}}$$

KT as an Optimization Algorithm

- Assume the function $F(\mathbf{x})$ is 1-Lipschitz ($G = 1$)
- $\mathbf{x}_t = \frac{\sum_{i=1}^{t-1} \mathbf{c}_i}{t} \text{Wealth}_{t-1} = \frac{-\sum_{i=1}^{t-1} \mathbf{g}_i}{t} (1 - \sum_{i=1}^{t-1} \langle \mathbf{g}_i, \mathbf{x}_i \rangle)$

Theorem (Informal)

KT betting in Hilbert spaces guarantees

$$F\left(\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t\right) - F(\mathbf{x}^*) \leq C \frac{\|\mathbf{x}^*\| \sqrt{\log(T)}}{\sqrt{T}}$$

Proof idea:

- “Worst” direction for gradient at time t is parallel to $\sum_{i=1}^{t-1} \mathbf{g}_i$
- “Worst” gradient have $\|\mathbf{g}_t\| = 1$

KT as an Optimization Algorithm

- Assume the function $F(\mathbf{x})$ is 1-Lipschitz ($G = 1$)
- $\mathbf{x}_t = \frac{\sum_{i=1}^{t-1} \mathbf{c}_i}{t} \text{Wealth}_{t-1} = \frac{-\sum_{i=1}^{t-1} \mathbf{g}_i}{t} (1 - \sum_{i=1}^{t-1} \langle \mathbf{g}_i, \mathbf{x}_i \rangle)$

Theorem (Informal)

KT betting in Hilbert spaces guarantees

$$F\left(\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t\right) - F(\mathbf{x}^*) \leq \min_{\alpha} C \frac{\left(\frac{\|\mathbf{x}^*\|^2}{\alpha} + \alpha\right) \sqrt{\log(T)}}{\sqrt{T}}$$

KT as an Optimization Algorithm

- Assume the function $F(\mathbf{x})$ is 1-Lipschitz ($G = 1$)
- $\mathbf{x}_t = \frac{\sum_{i=1}^{t-1} \mathbf{c}_i}{t} \text{Wealth}_{t-1} = \frac{-\sum_{i=1}^{t-1} \mathbf{g}_i}{t} (1 - \sum_{i=1}^{t-1} \langle \mathbf{g}_i, \mathbf{x}_i \rangle)$

Theorem (Informal)

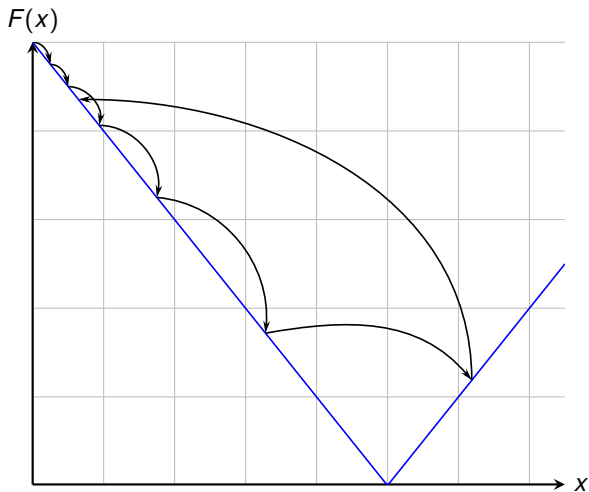
KT betting in Hilbert spaces guarantees

$$F\left(\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t\right) - F(\mathbf{x}^*) \leq \min_{\alpha} C \frac{\left(\frac{\|\mathbf{x}^*\|^2}{\alpha} + \alpha\right) \sqrt{\log(T)}}{\sqrt{T}}$$

- Compare it to the Gradient Descent guarantee with learning rate $\frac{\alpha}{\sqrt{t}}$

$$F\left(\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t\right) - F(\mathbf{x}^*) \leq C' \frac{\frac{\|\mathbf{x}^*\|^2}{\alpha} + \alpha}{\sqrt{T}}$$

How does the Betting Approach Work?



Constant Betting Case

- We said the optimal learning rate is $\eta_t = \frac{\|\mathbf{x}^*\|}{\sqrt{t}}$, when $\mathbf{x}_1 = \mathbf{0}$
- Can we approximate it with $\eta_t = \frac{\|\mathbf{x}_t\|}{\sqrt{t}}$ and use $\mathbf{x}_{t+1} = \mathbf{x}_t - \frac{\|\mathbf{x}_t\|}{\sqrt{t}} \mathbf{g}_t$?
[You et al., arXiv'17; You et al., ICPP'18; Bernstein et al., arXiv'20]
- ⚠ Such update rule would require $\mathbf{x}_1 \neq \mathbf{0}$

Constant Betting Case

- We said the optimal learning rate is $\eta_t = \frac{\|\mathbf{x}^*\|}{\sqrt{t}}$, when $\mathbf{x}_1 = \mathbf{0}$
- Can we approximate it with $\eta_t = \frac{\|\mathbf{x}_t\|}{\sqrt{T}}$ and use $\mathbf{x}_{t+1} = \mathbf{x}_t - \frac{\|\mathbf{x}_t\|}{\sqrt{T}} \mathbf{g}_t$?
[You et al., arXiv'17; You et al., ICPP'18; Bernstein et al., arXiv'20]
- ⚠ Such update rule would require $\mathbf{x}_1 \neq \mathbf{0}$
- Use a betting view: constant betting fraction of $\frac{1}{\sqrt{T}}$, i.e., $x_t = \frac{1}{\sqrt{T}} \text{Wealth}_{t-1}$
- ⚠ x_t is always positive
- $\text{Wealth}_t = \text{Wealth}_{t-1} + x_t c_t = \text{Wealth}_{t-1} (1 - \frac{1}{\sqrt{T}} g_t)$
- $x_{t+1} = \frac{1}{\sqrt{T}} \text{Wealth}_t = \frac{1}{\sqrt{T}} \text{Wealth}_{t-1} (1 - \frac{1}{\sqrt{T}} g_t) = x_t (1 - \frac{1}{\sqrt{T}} g_t) = x_t - \frac{x_t}{\sqrt{T}} g_t$

Constant Betting Case

- We said the optimal learning rate is $\eta_t = \frac{\|\mathbf{x}^*\|}{\sqrt{t}}$, when $\mathbf{x}_1 = \mathbf{0}$
- Can we approximate it with $\eta_t = \frac{\|\mathbf{x}_t\|}{\sqrt{T}}$ and use $\mathbf{x}_{t+1} = \mathbf{x}_t - \frac{\|\mathbf{x}_t\|}{\sqrt{T}} \mathbf{g}_t$?
[You et al., arXiv'17; You et al., ICPP'18; Bernstein et al., arXiv'20]

⚠ Such update rule would require $\mathbf{x}_1 \neq \mathbf{0}$

- Use a betting view: constant betting fraction of $\frac{1}{\sqrt{T}}$, i.e., $x_t = \frac{1}{\sqrt{T}} \text{Wealth}_{t-1}$

⚠ x_t is always positive

- $\text{Wealth}_t = \text{Wealth}_{t-1} + x_t c_t = \text{Wealth}_{t-1} (1 - \frac{1}{\sqrt{T}} g_t)$
- $x_{t+1} = \frac{1}{\sqrt{T}} \text{Wealth}_t = \frac{1}{\sqrt{T}} \text{Wealth}_{t-1} (1 - \frac{1}{\sqrt{T}} g_t) = x_t (1 - \frac{1}{\sqrt{T}} g_t) = x_t - \frac{x_t}{\sqrt{T}} g_t$
- Wealth lower bound: $\text{Wealth}_T = \prod_{t=1}^T (1 - \frac{1}{\sqrt{T}} g_t) \geq \exp(-\frac{1}{\sqrt{T}} \sum_{t=1}^T g_t - 1)$
- We get almost the optimal rate: $O(\frac{x^*}{\sqrt{T}} \ln(x^* \sqrt{T}))$ with $x^* \geq 0$

How to Go from 1d to \mathbb{R}^d ?

We have proposed 3 possible solutions:

- Under certain technical conditions, the algorithm can be transformed simply using inner products instead of multiplications:

$$\mathbf{x}_t = \frac{-\sum_{i=1}^{t-1} \mathbf{g}_i}{t} \left(1 - \sum_{i=1}^{t-1} \langle \mathbf{g}_i, \mathbf{x}_i \rangle \right).$$

How to Go from 1d to \mathbb{R}^d ?

We have proposed 3 possible solutions:

- Under certain technical conditions, the algorithm can be transformed simply using inner products instead of multiplications:

$$\mathbf{x}_t = \frac{-\sum_{i=1}^{t-1} \mathbf{g}_i}{t} \left(1 - \sum_{i=1}^{t-1} \langle \mathbf{g}_i, \mathbf{x}_i \rangle \right).$$

- Just using a coordinate-wise reduction: $x_{t,j} = \frac{-\sum_{i=1}^{t-1} g_{i,j}}{t} \left(1 - \sum_{i=1}^{t-1} g_{i,j} x_{i,j} \right).$

How to Go from 1d to \mathbb{R}^d ?

We have proposed 3 possible solutions:

- Under certain technical conditions, the algorithm can be transformed simply using inner products instead of multiplications:

$$\mathbf{x}_t = \frac{-\sum_{i=1}^{t-1} \mathbf{g}_i}{t} \left(1 - \sum_{i=1}^{t-1} \langle \mathbf{g}_i, \mathbf{x}_i \rangle \right).$$

- Just using a coordinate-wise reduction: $x_{t,j} = \frac{-\sum_{i=1}^{t-1} g_{i,j}}{t} \left(1 - \sum_{i=1}^{t-1} g_{i,j} x_{i,j} \right)$.
- Optimize the magnitude and direction of \mathbf{x}_t independently

Frustratingly Easy Magnitude and Direction Decomposition

Require: 1d algorithm \mathcal{A}_{1D} , unit ball $S \subset R^d$ algorithm \mathcal{A}_S

- 1: **for** $t = 1$ **to** T **do**
- 2: Get point $z_t \in \mathbb{R}$ from \mathcal{A}_{1D}
- 3: Get point $y_t \in S$ from \mathcal{A}_S
- 4: Play $x_t = z_t y_t \in R^d$
- 5: Receive gradient g_t
- 6: Send g_t as the t -th gradient to \mathcal{A}_S
- 7: Send $\langle g_t, y_t \rangle$ as the t -th gradient to \mathcal{A}_{1D}
- 8: **end for**

$$\begin{aligned} \text{Regret}_T(x^*) &= \sum_{t=1}^T \langle g_t, x_t - x^* \rangle = \sum_{t=1}^T \langle g_t, z_t y_t \rangle - \langle g_t, x^* \rangle \\ &= \sum_{t=1}^T \underbrace{(\langle g_t, y_t \rangle z_t - \langle g_t, y_t \rangle \|x^*\|)}_{\text{regret of } \mathcal{A}_{1D} \text{ at } \|x^*\| \in \mathbb{R}} + \sum_{t=1}^T \|x^*\| \underbrace{\left(\langle g_t, y_t \rangle - \left\langle g_t, \frac{x^*}{\|x^*\|} \right\rangle \right)}_{\text{regret of } \mathcal{A}_S \text{ at } \frac{x^*}{\|x^*\|} \in S} \\ &= \text{Regret}_T^{\mathcal{A}_{1D}}(\|x^*\|) + \|x^*\| \text{Regret}_T^{\mathcal{A}_S} \left(\frac{x^*}{\|x^*\|} \right) \end{aligned}$$

Summary of Part 2

- 1 Removing learning rates is indeed possible
- 2 You can make an exponential amount of money betting on a non-stochastic coin with a parameter-free algorithm
- 3 You can reduce optimization to betting on a coin
- 4 Hence, you can use a betting algorithm as a parameter-free optimization algorithm