

Parameter-free Online Optimization Part 3

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

Previously

- We saw that stochastic optimization can be solved via online linear optimization:

$$\mathbb{E} \left[\sum_{t=1}^T F(\mathbf{x}_t) - F(\mathbf{x}^*) \right] \leq \mathbb{E} \left[\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \right]$$
$$\mathbb{E}[\mathbf{g}_t] = \nabla F(\mathbf{x}_t)$$

- We saw that stochastic optimization can be solved via online linear optimization:

$$\mathbb{E} \left[\sum_{t=1}^T F(\mathbf{x}_t) - F(\mathbf{x}^*) \right] \leq \mathbb{E} \left[\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \right]$$
$$\mathbb{E}[\mathbf{g}_t] = \nabla F(\mathbf{x}_t)$$

- The optimally tuned gradient descent guarantees error $\frac{\|\mathbf{x}^* - \mathbf{x}_1\|_G}{\sqrt{T}}$ using learning rate $\eta = \frac{\|\mathbf{x}^* - \mathbf{x}_1\|}{G\sqrt{T}}$.

- We saw that stochastic optimization can be solved via online linear optimization:

$$\mathbb{E} \left[\sum_{t=1}^T F(\mathbf{x}_t) - F(\mathbf{x}^*) \right] \leq \mathbb{E} \left[\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \right]$$
$$\mathbb{E}[\mathbf{g}_t] = \nabla F(\mathbf{x}_t)$$

- The optimally tuned gradient descent guarantees error $\frac{\|\mathbf{x}^* - \mathbf{x}_1\|_G}{\sqrt{T}}$ using learning rate $\eta = \frac{\|\mathbf{x}^* - \mathbf{x}_1\|}{G\sqrt{T}}$.
- A parameter-free algorithm can obtain $\tilde{O} \left(\frac{\|\mathbf{x}^* - \mathbf{x}_1\|_G}{\sqrt{T}} \right)$. The logarithmic dependencies are tight.

- We saw that stochastic optimization can be solved via online linear optimization:

$$\mathbb{E} \left[\sum_{t=1}^T F(\mathbf{x}_t) - F(\mathbf{x}^*) \right] \leq \mathbb{E} \left[\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \right]$$
$$\mathbb{E}[\mathbf{g}_t] = \nabla F(\mathbf{x}_t)$$

- The optimally tuned gradient descent guarantees error $\frac{\|\mathbf{x}^* - \mathbf{x}_1\|_G}{\sqrt{T}}$ using learning rate $\eta = \frac{\|\mathbf{x}^* - \mathbf{x}_1\|}{G\sqrt{T}}$.
- A parameter-free algorithm can obtain $\tilde{O} \left(\frac{\|\mathbf{x}^* - \mathbf{x}_1\|_G}{\sqrt{T}} \right)$. The logarithmic dependencies are tight.
- What next?

In This Session:

Building algorithms that go beyond worst-case tuning.

In This Session:

Building algorithms that go beyond worst-case tuning.

- 1 How to build an “adaptive” parameter-free algorithm.

In This Session:

Building algorithms that go beyond worst-case tuning.

- 1 How to build an “adaptive” parameter-free algorithm.
- 2 Automatically improving when \mathbf{x}^* is sparse.

Building algorithms that go beyond worst-case tuning.

- 1 How to build an “adaptive” parameter-free algorithm.
- 2 Automatically improving when \mathbf{x}^* is sparse.
- 3 Automatically improving when $F(\mathbf{x})$ is smooth.

In This Session:

Building algorithms that go beyond worst-case tuning.

- 1 How to build an “adaptive” parameter-free algorithm.
- 2 Automatically improving when \mathbf{x}^* is sparse.
- 3 Automatically improving when $F(\mathbf{x})$ is smooth.
- 4 Automatically improving when $F(\mathbf{x})$ is strongly-convex.

Building algorithms that go beyond worst-case tuning.

- 1 How to build an “adaptive” parameter-free algorithm.
- 2 Automatically improving when \mathbf{x}^* is sparse.
- 3 Automatically improving when $F(\mathbf{x})$ is smooth.
- 4 Automatically improving when $F(\mathbf{x})$ is strongly-convex.
- 5 Dependencies on Lipschitz constants G .

Building algorithms that go beyond worst-case tuning.

- 1 How to build an “adaptive” parameter-free algorithm.
- 2 Automatically improving when \mathbf{x}^* is sparse.
- 3 Automatically improving when $F(\mathbf{x})$ is smooth.
- 4 Automatically improving when $F(\mathbf{x})$ is strongly-convex.
- 5 Dependencies on Lipschitz constants G .
 - We will assume $\mathbf{x}_1 = 0$ in this section to avoid writing $\mathbf{x}^* - \mathbf{x}_1$ too much.

Second-Order Regret Bound

We will build an online linear optimization algorithm that obtains:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \tilde{O} \left(\|\mathbf{x}^*\| \sqrt{\sum_{t=1}^T \|\mathbf{g}_t\|^2} \right)$$

This is the “AdaGrad version” of the KT-estimator regret bound $\tilde{O}(\|\mathbf{x}^*\| G\sqrt{T})$.

Second-Order Regret Bound

We will build an online linear optimization algorithm that obtains:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \tilde{O} \left(\|\mathbf{x}^*\| \sqrt{\sum_{t=1}^T \|\mathbf{g}_t\|^2} \right)$$

This is the “AdaGrad version” of the KT-estimator regret bound $\tilde{O}(\|\mathbf{x}^*\| G\sqrt{T})$.

⚠ This is called “second order” because of the $\|\mathbf{g}_t\|^2$ - it does NOT make any use of the Hessian.

Second-Order Regret Bound

We will build an online linear optimization algorithm that obtains:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \tilde{O} \left(\|\mathbf{x}^*\| \sqrt{\sum_{t=1}^T \|\mathbf{g}_t\|^2} \right)$$

This is the “AdaGrad version” of the KT-estimator regret bound $\tilde{O}(\|\mathbf{x}^*\| G \sqrt{T})$.

⚠ This is called “second order” because of the $\|\mathbf{g}_t\|^2$ - it does NOT make any use of the Hessian.

Let us see how to obtain this bound. For simplicity, we will assume $G = 1$.

Obtaining the Second-Order Bound (in 1d)

Recall the basic coin-betting strategy:

Obtaining the Second-Order Bound (in 1d)

Recall the basic coin-betting strategy:

- 1 Think of the gradient g_t as specifying (continuous) outcomes of “coin flips”.
At each iteration, we earn money $-g_t x_t$.

Obtaining the Second-Order Bound (in 1d)

Recall the basic coin-betting strategy:

- 1 Think of the gradient g_t as specifying (continuous) outcomes of “coin flips”.
At each iteration, we earn money $-g_t x_t$.
- 2 Keep track of the wealth: $\text{Wealth}_t = 1 - \sum_{i=1}^t g_i x_i$.

Obtaining the Second-Order Bound (in 1d)

Recall the basic coin-betting strategy:

- 1 Think of the gradient g_t as specifying (continuous) outcomes of “coin flips”.
At each iteration, we earn money $-g_t x_t$.
- 2 Keep track of the wealth: $\text{Wealth}_t = 1 - \sum_{i=1}^t g_i x_i$.
- 3 The output x_t corresponds to “betting” some money on the value of g_t .

Obtaining the Second-Order Bound (in 1d)

Recall the basic coin-betting strategy:

- 1 Think of the gradient g_t as specifying (continuous) outcomes of “coin flips”. At each iteration, we earn money $-g_t x_t$.
- 2 Keep track of the wealth: $\text{Wealth}_t = 1 - \sum_{i=1}^t g_i x_i$.
- 3 The output x_t corresponds to “betting” some money on the value of g_t .
- 4 If we can make $\text{Wealth}_T \geq H(\sum_{t=1}^T g_t)$ for any function H , then regret is bounded:

$$\sum_{t=1}^T g_t(x_t - x^*) \leq 1 + H^*(x^*)$$

where H^* is the Fenchel conjugate of H .

Obtaining the Second-Order Bound (in 1d)

Recall the basic coin-betting strategy:

- 1 Think of the gradient g_t as specifying (continuous) outcomes of “coin flips”. At each iteration, we earn money $-g_t x_t$.
- 2 Keep track of the wealth: $\text{Wealth}_t = 1 - \sum_{i=1}^t g_i x_i$.
- 3 The output x_t corresponds to “betting” some money on the value of g_t .
- 4 If we can make $\text{Wealth}_T \geq H(\sum_{t=1}^T g_t)$ for any function H , then regret is bounded:

$$\sum_{t=1}^T g_t(x_t - x^*) \leq 1 + H^*(x^*)$$

where H^* is the Fenchel conjugate of H .

- You are not allowed to borrow money: $x_t = \beta_t \text{Wealth}_{t-1}$, where $\beta_t \in [-1, 1]$ is called the “betting fraction”.

Obtaining the Second-Order Bound (in 1d)

Recall the basic coin-betting strategy:

- 1 Think of the gradient g_t as specifying (continuous) outcomes of “coin flips”. At each iteration, we earn money $-g_t x_t$.
- 2 Keep track of the wealth: $\text{Wealth}_t = 1 - \sum_{i=1}^t g_i x_i$.
- 3 The output x_t corresponds to “betting” some money on the value of g_t .
- 4 If we can make $\text{Wealth}_T \geq H(\sum_{t=1}^T g_t)$ for any function H , then regret is bounded:

$$\sum_{t=1}^T g_t(x_t - x^*) \leq 1 + H^*(x^*)$$

where H^* is the Fenchel conjugate of H .

- You are not allowed to borrow money: $x_t = \beta_t \text{Wealth}_{t-1}$, where $\beta_t \in [-1, 1]$ is called the “betting fraction”.
- This is actually no restriction: any algorithm that guarantees $\text{Regret}_T(0) = O(1)$ must be a coin-betting algorithm.

Obtaining the Second-Order Bound (in 1d)

Strategy:

- 1 Compute a near-optimal fixed β^* with foresight (similar to optimal Kelly bet).

Obtaining the Second-Order Bound (in 1d)

Strategy:

- 1 Compute a near-optimal fixed β^* with foresight (similar to optimal Kelly bet).
- 2 Observe that this fixed β^* yields the second-order bound.

Obtaining the Second-Order Bound (in 1d)

Strategy:

- 1 Compute a near-optimal fixed β^* with foresight (similar to optimal Kelly bet).
- 2 Observe that this fixed β^* yields the second-order bound.
- 3 Find a way to generate β_t on-the-fly that does nearly as well as β^* .

Calculating the Wealth

Since $x_t = \beta_t \text{Wealth}_{t-1}$, we have a recursion for Wealth_T

Calculating the Wealth

Since $x_t = \beta_t \text{Wealth}_{t-1}$, we have a recursion for Wealth_T

$$\text{Wealth}_T = 1 - \sum_{t=1}^T g_t x_t = \text{Wealth}_{T-1} - g_T x_T$$

Calculating the Wealth

Since $x_t = \beta_t \text{Wealth}_{t-1}$, we have a recursion for Wealth_T

$$\begin{aligned}\text{Wealth}_T &= 1 - \sum_{t=1}^T g_t x_t = \text{Wealth}_{T-1} - g_T x_T \\ &= \text{Wealth}_{T-1}(1 - \beta_T g_T) = \prod_{t=1}^T (1 - \beta_t g_t)\end{aligned}$$

Calculating the Wealth

Since $x_t = \beta_t \text{Wealth}_{t-1}$, we have a recursion for Wealth_T

$$\begin{aligned}\text{Wealth}_T &= 1 - \sum_{t=1}^T g_t x_t = \text{Wealth}_{T-1} - g_T x_T \\ &= \text{Wealth}_{T-1}(1 - \beta_T g_T) = \prod_{t=1}^T (1 - \beta_t g_t)\end{aligned}$$

Whenever you see a product, it's natural to take the logarithm:

$$\log(\text{Wealth}_T) = \sum_{t=1}^T \log(1 - \beta_t g_t)$$

A Good Fixed Betting Fraction

Let us guess $\beta^* = -\frac{\sum_{t=1}^T g_t}{\sum_{t=1}^T g_t^2 + 2|\sum_{t=1}^T g_t|}$. This yields $\beta^* \in [-1/2, 1/2]$, and so

A Good Fixed Betting Fraction

Let us guess $\beta^* = -\frac{\sum_{t=1}^T g_t}{\sum_{t=1}^T g_t^2 + 2|\sum_{t=1}^T g_t|}$. This yields $\beta^* \in [-1/2, 1/2]$, and so

$$\log(\text{Wealth}_T) = \sum_{t=1}^T \log(1 - \beta^* g_t)$$

apply a Taylor expansion of \log using $\beta^* g_t \in [-1/2, 1/2]$.

A Good Fixed Betting Fraction

Let us guess $\beta^* = -\frac{\sum_{t=1}^T g_t}{\sum_{t=1}^T g_t^2 + 2|\sum_{t=1}^T g_t|}$. This yields $\beta^* \in [-1/2, 1/2]$, and so

$$\log(\text{Wealth}_T) = \sum_{t=1}^T \log(1 - \beta^* g_t)$$

apply a Taylor expansion of log using $\beta^* g_t \in [-1/2, 1/2]$.

$$\geq \sum_{t=1}^T \left[-\beta^* g_t - (\beta^* g_t)^2 \right] \geq \frac{\left(\sum_{t=1}^T g_t \right)^2}{2 \sum_{t=1}^T g_t^2 + 4 \left| \sum_{t=1}^T g_t \right|}$$

A Good Fixed Betting Fraction

Let us guess $\beta^* = -\frac{\sum_{t=1}^T g_t}{\sum_{t=1}^T g_t^2 + 2|\sum_{t=1}^T g_t|}$. This yields $\beta^* \in [-1/2, 1/2]$, and so

$$\log(\text{Wealth}_T) = \sum_{t=1}^T \log(1 - \beta^* g_t)$$

apply a Taylor expansion of log using $\beta^* g_t \in [-1/2, 1/2]$.

$$\geq \sum_{t=1}^T \left[-\beta^* g_t - (\beta^* g_t)^2 \right] \geq \frac{\left(\sum_{t=1}^T g_t \right)^2}{2 \sum_{t=1}^T g_t^2 + 4 \left| \sum_{t=1}^T g_t \right|}$$

$$\text{Wealth}_T \geq \exp \left[\frac{\left(\sum_{t=1}^T g_t \right)^2}{2 \sum_{t=1}^T g_t^2 + 4 \left| \sum_{t=1}^T g_t \right|} \right]$$

A Good Fixed Betting Fraction

Let us guess $\beta^* = -\frac{\sum_{t=1}^T g_t}{\sum_{t=1}^T g_t^2 + 2|\sum_{t=1}^T g_t|}$. This yields $\beta^* \in [-1/2, 1/2]$, and so

$$\log(\text{Wealth}_T) = \sum_{t=1}^T \log(1 - \beta^* g_t)$$

apply a Taylor expansion of log using $\beta^* g_t \in [-1/2, 1/2]$.

$$\geq \sum_{t=1}^T \left[-\beta^* g_t - (\beta^* g_t)^2 \right] \geq \frac{\left(\sum_{t=1}^T g_t \right)^2}{2 \sum_{t=1}^T g_t^2 + 4 \left| \sum_{t=1}^T g_t \right|}$$

$$\text{Wealth}_T \geq \exp \left[\frac{\left(\sum_{t=1}^T g_t \right)^2}{2 \sum_{t=1}^T g_t^2 + 4 \left| \sum_{t=1}^T g_t \right|} \right] := H \left(\sum_{t=1}^T g_t \right)$$

A Good Fixed Betting Fraction

$$\text{Wealth}_T \geq \exp \left[\frac{\left(\sum_{t=1}^T g_t \right)^2}{2 \sum_{t=1}^T g_t^2 + 4 \left| \sum_{t=1}^T g_t \right|} \right] := H \left(\sum_{t=1}^T g_t \right)$$

A Good Fixed Betting Fraction

$$\text{Wealth}_T \geq \exp \left[\frac{\left(\sum_{t=1}^T g_t \right)^2}{2 \sum_{t=1}^T g_t^2 + 4 \left| \sum_{t=1}^T g_t \right|} \right] := H \left(\sum_{t=1}^T g_t \right)$$

$$\sum_{t=1}^T g_t (x_t - x^*) \leq 1 + H^*(x^*) = \tilde{O} \left(|x^*| \sqrt{\sum_{t=1}^T g_t^2} \right)$$

- 1 There is a fixed betting fraction β^* that would guarantee a second-order regret bound.
- 2 How can we learn about β^* on-the-fly?

Compute the gap between the log wealth we obtain vs the log wealth of β^* :

Compute the gap between the log wealth we obtain vs the log wealth of β^* :

$$\log(\text{Wealth}_T^*) - \log(\text{Wealth}_T) = \sum_{t=1}^T \log(1 - \beta^* g_t) - \sum_{t=1}^T \log(1 - \beta_t g_t)$$

Compute the gap between the log wealth we obtain vs the log wealth of β^* :

$$\log(\text{Wealth}_T^*) - \log(\text{Wealth}_T) = \sum_{t=1}^T \log(1 - \beta^* g_t) - \sum_{t=1}^T \log(1 - \beta_t g_t)$$

Define $\ell_t(\beta) = -\log(1 - \beta g_t)$:

$$:= \sum_{t=1}^T \ell_t(\beta_t) - \sum_{t=1}^T \ell_t(\beta^*)$$

Compute the gap between the log wealth we obtain vs the log wealth of β^* :

$$\log(\text{Wealth}_T^*) - \log(\text{Wealth}_T) = \sum_{t=1}^T \log(1 - \beta^* g_t) - \sum_{t=1}^T \log(1 - \beta_t g_t)$$

Define $\ell_t(\beta) = -\log(1 - \beta g_t)$:

$$:= \sum_{t=1}^T \ell_t(\beta_t) - \sum_{t=1}^T \ell_t(\beta^*)$$

ℓ_t is exp-concave! This means that choosing β_t by gradient descent on ℓ_t (e.g. via Online Newton Step [Hazan et al., MLJ'07]) yields:

$$\begin{aligned} &\leq \log(T) \\ \text{Wealth}_T &\geq \frac{\text{Wealth}_T^*}{T} \end{aligned}$$

Wrapping Up Second Order Bound

Now we have

$$\begin{aligned}\text{Wealth}_T &\geq \frac{\text{Wealth}_T^*}{T} \\ &\geq \frac{1}{T} \exp \left[\frac{\left(\sum_{t=1}^T g_t \right)^2}{2 \sum_{t=1}^T g_t^2 + 4 \left| \sum_{t=1}^T g_t \right|} \right] := H \left(\sum_{t=1}^T g_t \right) \\ \sum_{t=1}^T g_t (x_t - x^*) &\leq 1 + H^*(x^*) = \tilde{O} \left(|x^*| \sqrt{\sum_{t=1}^T g_t^2} \right)\end{aligned}$$

Wrapping Up Second Order Bound

Now we have

$$\begin{aligned}\text{Wealth}_T &\geq \frac{\text{Wealth}_T^*}{T} \\ &\geq \frac{1}{T} \exp \left[\frac{\left(\sum_{t=1}^T g_t \right)^2}{2 \sum_{t=1}^T g_t^2 + 4 \left| \sum_{t=1}^T g_t \right|} \right] := H \left(\sum_{t=1}^T g_t \right) \\ \sum_{t=1}^T g_t (x_t - x^*) &\leq 1 + H^*(x^*) = \tilde{O} \left(|x^*| \sqrt{\sum_{t=1}^T g_t^2} \right)\end{aligned}$$

The $\frac{1}{T}$ only changes a $\log(1 + |x^*| \sqrt{T})$ to $\log(1 + |x^*| T)$.

Full Second-Order Bound Algorithm

- 1: $\beta_1 = 0$, $\text{Wealth}_0 = 1$
- 2: **for** $t = 1$ **to** T **do**
- 3: Play $x_t = \beta_t \text{Wealth}_{t-1}$
- 4: Get gradient g_t , define $\ell_t(\beta) = -\log(1 - \beta g_t)$
- 5: Compute $z_t = \ell'_t(\beta_t) = \frac{g_t}{1 - \beta g_t}$
- 6: Set $\beta_{t+1} = \text{clip}\left(\beta_t - \frac{z_t}{1 + \sum_{i=1}^t z_i^2}, -\frac{1}{2}, \frac{1}{2}\right)$
 This is the ONS update
- 7: Set $\text{Wealth}_t = \text{Wealth}_{t-1} - g_t x_t$
- 8: **end for**

Theorem

The above algorithm guarantees

$$\sum_{t=1}^T g_t(x_t - x^*) \leq \tilde{O}\left(|x^*| \sqrt{\sum_{t=1}^T g_t^2 \log(|x^*| T)}\right)$$

Now we are armed with an second-order 1-d regret bound. We'll show how it can be easily used to accomplish some interesting tasks:

Now we are armed with an second-order 1-d regret bound. We'll show how it can be easily used to accomplish some interesting tasks:

- 1 A bound that depends on the sparsity of \mathbf{x}^* .

Now we are armed with an second-order 1-d regret bound. We'll show how it can be easily used to accomplish some interesting tasks:

- 1 A bound that depends on the sparsity of \mathbf{x}^* .
- 2 A bound that automatically adapts to the smoothness and variance parameters.

Now we are armed with an second-order 1-d regret bound. We'll show how it can be easily used to accomplish some interesting tasks:

- 1 A bound that depends on the sparsity of \mathbf{x}^* .
- 2 A bound that automatically adapts to the smoothness and variance parameters.
- 3 A bound that automatically adapts to strong-convexity parameters.

Sparsity via Per-Coordinate Updates

Remember one way to make a d -dimensional algorithm from a 1-dimensional algorithm, ala AdaGrad:

$$\begin{aligned} \sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle &= \sum_{i=1}^d \underbrace{\sum_{t=1}^T g_{t,i} (x_{t,i} - x_i^*)}_{\text{Set } \mathbf{x}_{t,i} \text{ via a 1-d learner.}} \\ &\leq \underbrace{\tilde{O} \left(\sum_{i=1}^d |x_i^*| \sqrt{\sum_{t=1}^T g_{t,i}^2} \right)}_{\text{The } \tilde{O} \text{ now hides a } \log(d) \text{ factor.}} \end{aligned}$$

Sparsity via Per-Coordinate Updates

Remember one way to make a d -dimensional algorithm from a 1-dimensional algorithm, ala AdaGrad:

$$\begin{aligned} \sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle &= \sum_{i=1}^d \underbrace{\sum_{t=1}^T g_{t,i} (x_{t,i} - x_i^*)}_{\text{Set } \mathbf{x}_{t,i} \text{ via a 1-d learner.}} \\ &\leq \underbrace{\tilde{O} \left(\sum_{i=1}^d |x_i^*| \sqrt{\sum_{t=1}^T g_{t,i}^2} \right)}_{\text{The } \tilde{O} \text{ now hides a } \log(d) \text{ factor.}} \end{aligned}$$

Contrast with the (optimally tuned) AdaGrad guarantee:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq O \left(\max_i |x_i^*| \sum_{i=1}^d \sqrt{\sum_{t=1}^T g_{t,i}^2} \right)$$

Sparsity via Per-Coordinate Updates

What happens if some coordinates x_i^* are 0?

- Suppose $\|\mathbf{x}^*\|_\infty \leq 1$, $\|\mathbf{g}_t\|_\infty \leq 1$, and \mathbf{x}^* is S -sparse. Then:

Sparsity via Per-Coordinate Updates

What happens if some coordinates x_i^* are 0?

- Suppose $\|\mathbf{x}^*\|_\infty \leq 1$, $\|\mathbf{g}_t\|_\infty \leq 1$, and \mathbf{x}^* is S -sparse. Then:

$$\begin{aligned}\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle &\leq \tilde{O} \left(\sum_{i=1}^d |x_i^*| \sqrt{\sum_{t=1}^T g_{t,i}^2} \right) \\ &\leq \tilde{O} \left(\|\mathbf{x}^*\|_1 \sqrt{\sum_{t=1}^T \|\mathbf{g}_t\|_\infty^2} \right) \\ &\leq S\sqrt{T}\end{aligned}$$

More Properties of the Per-Coordinate Update

Per-Coordinate updates also achieve the L_2 -norm bound:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \tilde{O} \left(\sum_{i=1}^d |x_i^*| \sqrt{\sum_{t=1}^T g_{t,i}^2} \right)$$

apply Cauchy-Schwarz $\leq \tilde{O} \left(\|\mathbf{x}^*\| \sqrt{\sum_{t=1}^T \|\mathbf{g}_t\|^2} \right)$

Both of these properties leverage the automatic adaptivity to each individual x_i^* .

More Properties of the Per-Coordinate Update

Per-Coordinate updates also achieve the L_2 -norm bound:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \tilde{O} \left(\sum_{i=1}^d |x_i^*| \sqrt{\sum_{t=1}^T g_{t,i}^2} \right)$$

apply Cauchy-Schwarz $\leq \tilde{O} \left(\|\mathbf{x}^*\| \sqrt{\sum_{t=1}^T \|\mathbf{g}_t\|^2} \right)$

Both of these properties leverage the automatic adaptivity to each individual x_i^* .

⚠ This is not the same as identifying which coordinates of \mathbf{x}^* are zero.

Parameter-Free Experts Algorithms

- Learning with expert advice setting is a special case of online linear optimization in which both \mathbf{x}_t and \mathbf{x}^* must lie in the probability simplex in \mathbb{R}^d .
- There are many parameter-free algorithms for this special case, starting with NormalHedge [Chaudhuri et al., NeurIPS'09], and including [Gaillard et al., COLT'14; Koolen&van Erven, COLT'15; Luo&Schapire, COLT'15; Foster et al., NeurIPS'15; Harvey et al., arXiv'20].
- In this special setting, one can hope for interesting different bounds. For example, the Squint algorithm can obtain:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \tilde{O} \left(\sqrt{\sum_{t=1}^T \sum_{i=1}^d |x_i^*| (\langle \mathbf{g}_t, \mathbf{x}_t \rangle - g_{t,i})^2} \right)$$

Parameter-Free Experts Algorithms

- Learning with expert advice setting is a special case of online linear optimization in which both \mathbf{x}_t and \mathbf{x}^* must lie in the probability simplex in \mathbb{R}^d .
- There are many parameter-free algorithms for this special case, starting with NormalHedge [Chaudhuri et al., NeurIPS'09], and including [Gaillard et al., COLT'14; Koolen&van Erven, COLT'15; Luo&Schapire, COLT'15; Foster et al., NeurIPS'15; Harvey et al., arXiv'20].
- In this special setting, one can hope for interesting different bounds. For example, the Squint algorithm can obtain:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \tilde{O} \left(\sqrt{\sum_{t=1}^T \sum_{i=1}^d |x_i^*| (\langle \mathbf{g}_t, \mathbf{x}_t \rangle - g_{t,i})^2} \right)$$

which implies:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \tilde{O} \left(\sum_{i=1}^d |x_i^*| \sqrt{\sum_{t=1}^T |g_{t,i}|} \right)$$

Parameter-Free Experts Algorithms

- Learning with expert advice setting is a special case of online linear optimization in which both \mathbf{x}_t and \mathbf{x}^* must lie in the probability simplex in \mathbb{R}^d .
- There are many parameter-free algorithms for this special case, starting with NormalHedge [Chaudhuri et al., NeurIPS'09], and including [Gaillard et al., COLT'14; Koolen&van Erven, COLT'15; Luo&Schapire, COLT'15; Foster et al., NeurIPS'15; Harvey et al., arXiv'20].
- In this special setting, one can hope for interesting different bounds. For example, the Squint algorithm can obtain:

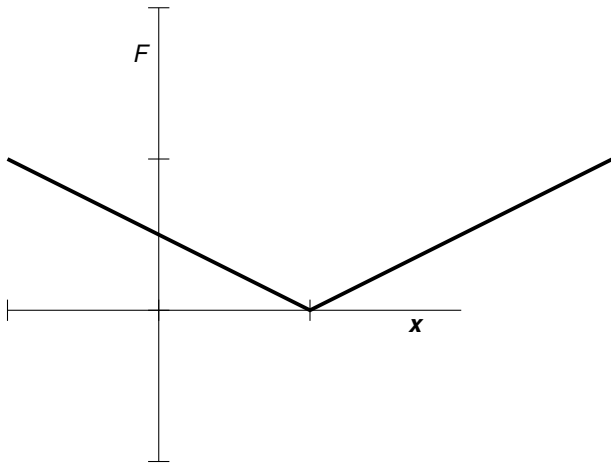
$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \tilde{O} \left(\sqrt{\sum_{t=1}^T \sum_{i=1}^d |x_i^*| (\langle \mathbf{g}_t, \mathbf{x}_t \rangle - g_{t,i})^2} \right)$$

which implies:

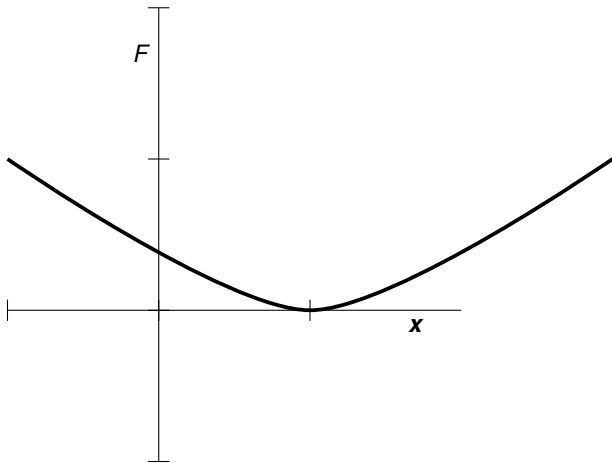
$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \tilde{O} \left(\sum_{i=1}^d |x_i^*| \sqrt{\sum_{t=1}^T |g_{t,i}|} \right)$$

- This does not have the $g_{t,i}^2$, but works in the simplex.

Adapting to Smoothness



Adapting to Smoothness



Adapting to Smoothness

Suppose our original objective $F(\mathbf{x})$ is L -smooth:

$$\nabla^2 F(\mathbf{x}) \preceq LI$$

We don't know what L is. Suppose also that for all t ,

$$\text{Var}[\mathbf{g}_t] \leq \sigma^2$$

In this case, the optimal learning rate for gradient descent depends on L and σ , and achieves:

$$\mathbb{E}[F(\mathbf{x}_{\text{SGD}}) - F(\mathbf{x}^*)] \leq \frac{L\|\mathbf{x}^*\|^2}{T} + \frac{\sigma\|\mathbf{x}^*\|}{\sqrt{T}}$$

- The second-order parameter-free bound automatically obtains this result up to a log factor!
- To see how, we need to look inside the online-to-batch conversion, and look beyond the linear approximation to F .

- The second-order parameter-free bound automatically obtains this result up to a log factor!
- To see how, we need to look inside the online-to-batch conversion, and look beyond the linear approximation to F .
- Specifically, we need the following simple consequence of smoothness (see [Srebro et al., NeurIPS'10; Levy, NeurIPS'17] for related applications):

$$\mathbb{E}[\|\mathbf{g}_t\|^2] \leq L(F(\mathbf{x}_t) - F(\mathbf{x}^*)) + \sigma^2$$

Adapting to Smoothness

In the online-to-batch conversion, we argued:

$$\mathbb{E} \left[\sum_{t=1}^T F(\mathbf{x}_t) - F(\mathbf{x}^*) \right] \leq \mathbb{E} \left[\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \right] \text{ (convexity)}$$

Adapting to Smoothness

In the online-to-batch conversion, we argued:

$$\mathbb{E} \left[\sum_{t=1}^T F(\mathbf{x}_t) - F(\mathbf{x}^*) \right] \leq \mathbb{E} \left[\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \right] \quad (\text{convexity})$$

use regret bound $\leq \mathbb{E} \left[\tilde{O} \left(\|\mathbf{x}^*\| \sqrt{\sum_{t=1}^T \|\mathbf{g}_t\|^2} \right) \right]$

Adapting to Smoothness

In the online-to-batch conversion, we argued:

$$\mathbb{E} \left[\sum_{t=1}^T F(\mathbf{x}_t) - F(\mathbf{x}^*) \right] \leq \mathbb{E} \left[\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \right] \quad (\text{convexity})$$

$$\text{use regret bound} \leq \mathbb{E} \left[\tilde{O} \left(\|\mathbf{x}^*\| \sqrt{\sum_{t=1}^T \|\mathbf{g}_t\|^2} \right) \right]$$

$$\text{use Jensen} \leq \tilde{O} \left(\|\mathbf{x}^*\| \sqrt{\mathbb{E} \left[\sum_{t=1}^T \|\mathbf{g}_t\|^2 \right]} \right)$$

Adapting to Smoothness

In the online-to-batch conversion, we argued:

$$\mathbb{E} \left[\sum_{t=1}^T F(\mathbf{x}_t) - F(\mathbf{x}^*) \right] \leq \mathbb{E} \left[\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \right] \quad (\text{convexity})$$

$$\text{use regret bound} \leq \mathbb{E} \left[\tilde{O} \left(\|\mathbf{x}^*\| \sqrt{\sum_{t=1}^T \|\mathbf{g}_t\|^2} \right) \right]$$

$$\text{use Jensen} \leq \tilde{O} \left(\|\mathbf{x}^*\| \sqrt{\mathbb{E} \left[\sum_{t=1}^T \|\mathbf{g}_t\|^2 \right]} \right)$$

$$\text{use smoothness} \leq \tilde{O} \left(\|\mathbf{x}^*\| \sqrt{L \mathbb{E} \left[\sum_{t=1}^T F(\mathbf{x}_t) - F(\mathbf{x}^*) \right] + T\sigma^2} \right)$$

Adapting to Smoothness

In the online-to-batch conversion, we argued:

$$\mathbb{E} \left[\sum_{t=1}^T F(\mathbf{x}_t) - F(\mathbf{x}^*) \right] \leq \mathbb{E} \left[\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \right] \quad (\text{convexity})$$

$$\text{use regret bound} \leq \mathbb{E} \left[\tilde{O} \left(\|\mathbf{x}^*\| \sqrt{\sum_{t=1}^T \|\mathbf{g}_t\|^2} \right) \right]$$

$$\text{use Jensen} \leq \tilde{O} \left(\|\mathbf{x}^*\| \sqrt{\mathbb{E} \left[\sum_{t=1}^T \|\mathbf{g}_t\|^2 \right]} \right)$$

$$\text{use smoothness} \leq \tilde{O} \left(\|\mathbf{x}^*\| \sqrt{L \mathbb{E} \left[\sum_{t=1}^T F(\mathbf{x}_t) - F(\mathbf{x}^*) \right] + T\sigma^2} \right)$$

$$\mathbb{E} \left[\sum_{t=1}^T F(\mathbf{x}_t) - F(\mathbf{x}^*) \right] \leq \tilde{O} \left(L\|\mathbf{x}^*\|^2 + \sigma\|\mathbf{x}^*\|\sqrt{T} \right) \quad (\text{quadratic formula})$$

Adapting to Smoothness

In the online-to-batch conversion, we argued:

$$\mathbb{E} \left[\sum_{t=1}^T F(\mathbf{x}_t) - F(\mathbf{x}^*) \right] \leq \mathbb{E} \left[\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \right] \quad (\text{convexity})$$

$$\text{use regret bound} \leq \mathbb{E} \left[\tilde{O} \left(\|\mathbf{x}^*\| \sqrt{\sum_{t=1}^T \|\mathbf{g}_t\|^2} \right) \right]$$

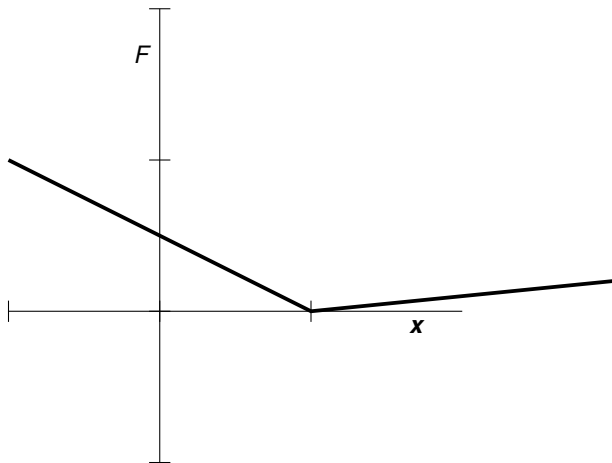
$$\text{use Jensen} \leq \tilde{O} \left(\|\mathbf{x}^*\| \sqrt{\mathbb{E} \left[\sum_{t=1}^T \|\mathbf{g}_t\|^2 \right]} \right)$$

$$\text{use smoothness} \leq \tilde{O} \left(\|\mathbf{x}^*\| \sqrt{L \mathbb{E} \left[\sum_{t=1}^T F(\mathbf{x}_t) - F(\mathbf{x}^*) \right] + T\sigma^2} \right)$$

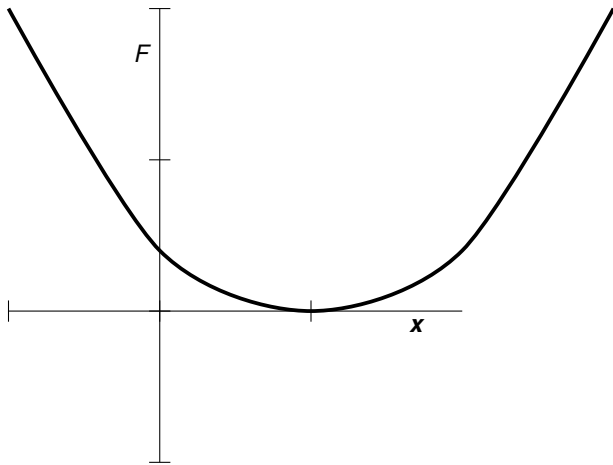
$$\mathbb{E} \left[\sum_{t=1}^T F(\mathbf{x}_t) - F(\mathbf{x}^*) \right] \leq \tilde{O} \left(L\|\mathbf{x}^*\|^2 + \sigma\|\mathbf{x}^*\|\sqrt{T} \right) \quad (\text{quadratic formula})$$

$$\mathbb{E} \left[F \left(\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \right) \right] - F(\mathbf{x}^*) \leq \tilde{O} \left(\frac{L\|\mathbf{x}^*\|^2}{T} + \frac{\sigma\|\mathbf{x}^*\|}{\sqrt{T}} \right) \quad (\text{Jensen})$$

Adapting to Strong Convexity



Adapting to Strong Convexity



Adapting to Strong-Convexity

- The strongly-convex losses are an important special class of optimization problems:

$$\nabla^2 F(\mathbf{x}) \succeq \mu I$$

- For strongly-convex losses, gradient descent with learning rate $\frac{1}{\mu t}$ guarantees [Hazan et al., MLJ'07]:

$$\mathbb{E}[F(\mathbf{x}_{\text{SGD}}) - F(\mathbf{x}^*)] \leq \frac{\log(T)}{\mu T}$$

This requires knowledge of the parameter μ . Can we do without this?

A Regret Bound for Strong-Convexity

- MetaGrad [Koolan&van Erven, NeurIPS'16] and Maler [Wang et al., UAI'19] are online linear optimization algorithms that guarantee:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \tilde{O} \left(\sqrt{\sum_{t=1}^T \|\mathbf{x}_t - \mathbf{x}^*\|^2 \|\mathbf{g}_t\|^2} \right)$$

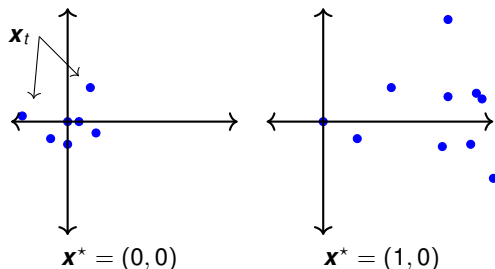
- Amazingly, this automatically implies adaptivity to strong-convexity!
- We'll look at a simple trick for obtaining this style of bound.

Recall:

- We obtain error $\tilde{O}\left(\frac{\|\mathbf{x}^*\|_G}{\sqrt{T}}\right)$ (ignoring logs and constants).
- Intuitively, when \mathbf{x}^* is far from the starting point, we have more space to explore, so we expect to \mathbf{x}_t to move a lot..

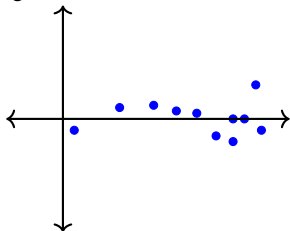
When \mathbf{x}^* is far from the starting point, \mathbf{x}_t should move a lot more.

When \mathbf{x}^* is far from the starting point, \mathbf{x}_t should move a lot more.



- Get much “tighter” distribution of \mathbf{x}_t s when \mathbf{x}^* is small.

Can we have a tight distribution even when \mathbf{x}^* is not small?



- We can get suboptimality

$$\mathbb{E} \left[F \left(\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \right) - F(\mathbf{x}^*) \right] = \tilde{O} \left(\frac{\|\mathbf{x}^*\|_G}{\sqrt{T}} \right)$$

- We can get suboptimality

$$\mathbb{E} \left[F \left(\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \right) - F(\mathbf{x}^*) \right] = \tilde{O} \left(\frac{\|\mathbf{x}^*\|G}{\sqrt{T}} \right)$$

- If we shift all outputs \mathbf{x}_t by a constant \mathbf{v} , then we get

$$\mathbb{E} \left[F \left(\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \right) - F(\mathbf{x}^*) \right] = \tilde{O} \left(\frac{\|\mathbf{x}^* - \mathbf{v}\|G}{\sqrt{T}} \right)$$

- We can get suboptimality

$$\mathbb{E} \left[F \left(\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \right) - F(\mathbf{x}^*) \right] = \tilde{O} \left(\frac{\|\mathbf{x}^*\|G}{\sqrt{T}} \right)$$

- If we shift all outputs \mathbf{x}_t by a constant \mathbf{v} , then we get

$$\mathbb{E} \left[F \left(\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \right) - F(\mathbf{x}^*) \right] = \tilde{O} \left(\frac{\|\mathbf{x}^* - \mathbf{v}\|G}{\sqrt{T}} \right)$$

- If \mathbf{v} is close to \mathbf{x}^* , then this is much better.

- We can get suboptimality

$$\mathbb{E} \left[F \left(\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \right) - F(\mathbf{x}^*) \right] = \tilde{O} \left(\frac{\|\mathbf{x}^*\|G}{\sqrt{T}} \right)$$

- If we shift all outputs \mathbf{x}_t by a constant \mathbf{v} , then we get

$$\mathbb{E} \left[F \left(\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \right) - F(\mathbf{x}^*) \right] = \tilde{O} \left(\frac{\|\mathbf{x}^* - \mathbf{v}\|G}{\sqrt{T}} \right)$$

- If \mathbf{v} is close to \mathbf{x}^* , then this is much better.
- So how can we find a closer \mathbf{v} ?

Using Strong Convexity

We need to move beyond the naive linear approximation

$F(\mathbf{x}_t) - F(\mathbf{x}^*) \leq \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle$. Strongly-convex functions also satisfy:

$$F(\mathbf{x}) - F(\mathbf{x}^*) \geq \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}^*\|^2$$

This means that points with low error are close to \mathbf{x}^* .

Adapting to Strong Convexity

- Small function value means that we are close to \mathbf{x}^* .

Adapting to Strong Convexity

- Small function value means that we are close to \mathbf{x}^* .
- The “average iterate” $\frac{\sum_{t=1}^T \mathbf{x}_t}{T}$ has small function value, so it must be close to \mathbf{x}^* !

Adapting to Strong Convexity

- Small function value means that we are close to \mathbf{x}^* .
- The “average iterate” $\frac{\sum_{t=1}^T \mathbf{x}_t}{T}$ has small function value, so it must be close to \mathbf{x}^* ! We should try using a running average for \mathbf{v} :

Adapting to Strong Convexity

- Small function value means that we are close to \mathbf{x}^* .
- The “average iterate” $\frac{\sum_{t=1}^T \mathbf{x}_t}{T}$ has small function value, so it must be close to \mathbf{x}^* ! We should try using a running average for \mathbf{v} :

$\hat{\mathbf{x}}_{t+1}$ = Parameter-free Output

$$\mathbf{x}_{t+1} = \hat{\mathbf{x}}_{t+1} + \frac{\sum_{i=1}^t \|\mathbf{g}_i\|^2 \mathbf{x}_i}{\underbrace{\sum_{i=1}^t \|\mathbf{g}_i\|^2}_{\text{Can use } \sum_{i=1}^t \mathbf{x}_i / t \text{ too.}}}$$

Adapting to Strong Convexity

- Small function value means that we are close to \mathbf{x}^* .
- The “average iterate” $\frac{\sum_{t=1}^T \mathbf{x}_t}{T}$ has small function value, so it must be close to \mathbf{x}^* ! We should try using a running average for v :

$\hat{\mathbf{x}}_{t+1}$ = Parameter-free Output

$$\mathbf{x}_{t+1} = \hat{\mathbf{x}}_{t+1} + \underbrace{\frac{\sum_{i=1}^t \|\mathbf{g}_i\|^2 \mathbf{x}_i}{\sum_{i=1}^t \|\mathbf{g}_i\|^2}}_{\text{Can use } \sum_{i=1}^t \mathbf{x}_i / t \text{ too.}}$$

- Via some careful algebra, we get:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \tilde{O} \left(\sqrt{\sum_{t=1}^T \|\mathbf{x}_t - \mathbf{x}^*\|^2 \|\mathbf{g}_t\|^2} \right)$$

This implies $\tilde{O} \left(\frac{1}{\mu T} \right)$ error in strongly-convex problems.

Adapting to Strong Convexity

- Small function value means that we are close to \mathbf{x}^* .
- The “average iterate” $\frac{\sum_{t=1}^T \mathbf{x}_t}{T}$ has small function value, so it must be close to \mathbf{x}^* ! We should try using a running average for v :

$\hat{\mathbf{x}}_{t+1}$ = Parameter-free Output

$$\mathbf{x}_{t+1} = \hat{\mathbf{x}}_{t+1} + \underbrace{\frac{\sum_{i=1}^t \|\mathbf{g}_i\|^2 \mathbf{x}_i}{\sum_{i=1}^t \|\mathbf{g}_i\|^2}}_{\text{Can use } \sum_{i=1}^t \mathbf{x}_i / t \text{ too.}}$$

- Via some careful algebra, we get:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \tilde{O} \left(\sqrt{\sum_{t=1}^T \|\mathbf{x}_t - \mathbf{x}^*\|^2 \|\mathbf{g}_t\|^2} \right)$$

This implies $\tilde{O} \left(\frac{1}{\mu T} \right)$ error in strongly-convex problems.

- Adding the average works for any parameter-free algorithm.

Dealing with Unknown G

- We have been assuming $G = 1$. In practice, one can usually guess a good upper bound for G by just looking at the first few gradients.

Dealing with Unknown G

- We have been assuming $G = 1$. In practice, one can usually guess a good upper bound for G by just looking at the first few gradients.
- The guess G can be wildly wrong before anything bad happens:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \tilde{O} \left(\|\mathbf{x}^*\| \sqrt{G^2 + \sum_{t=1}^T \|\mathbf{g}_t\|^2} \right)$$

Dealing with Unknown G

- We have been assuming $G = 1$. In practice, one can usually guess a good upper bound for G by just looking at the first few gradients.
- The guess G can be wildly wrong before anything bad happens:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \tilde{O} \left(\|\mathbf{x}^*\| \sqrt{G^2 + \sum_{t=1}^T \|\mathbf{g}_t\|^2} \right)$$

- It is “safe” for G to over-estimate by a factor of $O(\sqrt{T})$.

Dealing with Unknown G

- We have been assuming $G = 1$. In practice, one can usually guess a good upper bound for G by just looking at the first few gradients.
- The guess G can be wildly wrong before anything bad happens:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \tilde{O} \left(\|\mathbf{x}^*\| \sqrt{G^2 + \sum_{t=1}^T \|\mathbf{g}_t\|^2} \right)$$

- It is “safe” for G to over-estimate by a factor of $O(\sqrt{T})$.
- In contrast, the classical adaptive gradient approach stipulates a known bound $B \geq \|\mathbf{x}^*\|$ to guarantee:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq O \left(B \sqrt{\sum_{t=1}^T \|\mathbf{g}_t\|^2} \right)$$

Dealing with Unknown G

- We have been assuming $G = 1$. In practice, one can usually guess a good upper bound for G by just looking at the first few gradients.
- The guess G can be wildly wrong before anything bad happens:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \tilde{O} \left(\|\mathbf{x}^*\| \sqrt{G^2 + \sum_{t=1}^T \|\mathbf{g}_t\|^2} \right)$$

- It is “safe” for G to over-estimate by a factor of $O(\sqrt{T})$.
- In contrast, the classical adaptive gradient approach stipulates a known bound $B \geq \|\mathbf{x}^*\|$ to guarantee:

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq O \left(B \sqrt{\sum_{t=1}^T \|\mathbf{g}_t\|^2} \right)$$

- This pays linearly for any over-estimate in B .

- What about the true worst-case bound?

- What about the true worst-case bound?
- If we do not know G , an adversary can force regret $\tilde{\Omega}(\|\mathbf{x}^*\|G\sqrt{T} + G\|\mathbf{x}^*\|^3)$ [Mhammedi&Koolen, COLT'20].

- What about the true worst-case bound?
- If we do not know G , an adversary can force regret $\tilde{\Omega}(\|\mathbf{x}^*\|G\sqrt{T} + G\|\mathbf{x}^*\|^3)$ [Mhammedi&Koolen, COLT'20].
 - Note that the extra $\|\mathbf{x}^*\|^3$ term does not depend on T .

- What about the true worst-case bound?
- If we do not know G , an adversary can force regret $\tilde{\Omega}(\|\mathbf{x}^*\|G\sqrt{T} + G\|\mathbf{x}^*\|^3)$ [Mhammedi&Koolen, COLT'20].
 - Note that the extra $\|\mathbf{x}^*\|^3$ term does not depend on T .
 - It is unclear how “real” this lower bound is in real stochastic environments.

- What about the true worst-case bound?
- If we do not know G , an adversary can force regret $\tilde{\Omega}(\|\mathbf{x}^*\|G\sqrt{T} + G\|\mathbf{x}^*\|^3)$ [Mhammedi&Koolen, COLT'20].
 - Note that the extra $\|\mathbf{x}^*\|^3$ term does not depend on T .
 - It is unclear how “real” this lower bound is in real stochastic environments.
- The bound is tight. It can be obtained “lying” to the algorithm by feeding \mathbf{g}_t clipped to have norm at most $\max_{i < t} \|\mathbf{g}_i\|$ [Cutkosky, COLT'19; Mhammedi&Koolen, COLT'20].

Summary of Part 3

Once you know how to make parameter-free algorithms, it is surprisingly easy to build extra results.

- 1 Adapt to sparse \mathbf{x}^* .
- 2 Adapt to smoothness.
- 3 Adapt to strong-convexity.
- 4 Incorporate unknown G .