

Lecture 16 Online Newton Method

Recap: Follow the Regularized Leader

$$R(x) = \min_{x \in K} \left\langle x, \sum_{t=1}^T \nabla_t \right\rangle + R(x)$$

$\nabla_t = \nabla f_t(x_t)$

Regret $\leq D_R G_R \sqrt{2T}$

$$D_R = \int \max_{x, y \in K} (R(x) - R(y))$$

$$G_R \geq \|\nabla_t\|_t^* \quad \forall t$$

$$\|\cdot\|_t = \|\cdot\|_{\nabla R(z)}$$

Recall: Key Lemma (Stability)

$$\text{Regret} \leq \sum_{t=1}^T \nabla_t^T (x_t - x_{t+1}) + \frac{1}{\eta} D_R^2$$

Long argument: $\nabla_t^T (x_t - x_{t+1}) \leq 2\eta \|\nabla_t\|_t^*$

1. we bounded

$$\sum_t \nabla_t^T x_t - \sum_t \nabla_t^T u$$

$$\text{Regret} \leq \sum_{t=1}^T f_t(x_t) - \sum_{t=1}^T f_t(x^*)$$

why is this ok?

$$\rightarrow \sum_{t=1}^T f_t(x_t) - f_t(u) \leq \sum_{t=1}^T \nabla_t^T (x_t - u)$$

Convex f_t

2. $R(x) = x^T \log x = \sum_{i=1}^n x_i \log x_i$ (negative entropy)

$K = \text{Probability Simplex} = \Delta_n = \{x \mid x \geq 0, \mathbf{1}^T x = 1\}$

$G_R = 1, D_R \leq \sqrt{\log n}$ $\sqrt{T \log n}$

Recovers the Hedge algorithm / exponentiated gradients

Alternative View from Strong Convexity

$$\nabla_t^T (x_t - x_{t+1}) \leq \|\nabla_t\|_t \|x_t - x_{t+1}\|_t = \|x_t - x_{t+1}\|_t$$

$$\sum \nabla_t^T (x_t - x_{t+1}) \leq \eta T$$

$R(x) = x^T \log x$ is 1-strongly convex in $\Delta_n \in \text{AW4}$
for all $x \in \Delta_n$

$\|x_{t+1} - x_t\|_1 \leq ?$

$$F(x) = \underbrace{\langle \nabla F(x), x \rangle} + R(x)$$

$$F(x_{t+1}) \geq F(x_t) + \langle \nabla F(x_t), x_{t+1} - x_t \rangle + \frac{1}{2} \|x_{t+1} - x_t\|_1^2$$

$$F(x_t) \geq F(x_{t+1}) + \langle \nabla F(x_{t+1}), x_t - x_{t+1} \rangle - \frac{1}{2} \|x_t - x_{t+1}\|_1^2$$

$$0 \geq \langle x_{t+1} - x_t, \nabla F(x_t) - \nabla F(x_{t+1}) \rangle + \|x_t - x_{t+1}\|_1^2$$

$$\|x_t - x_{t+1}\|_1^2 \leq \|x_{t+1} - x_t\|_1 \|\nabla F(x_t) - \nabla F(x_{t+1})\|_\infty$$

$$\|x_t - x_{t+1}\|_1 \leq \eta$$

Universal Portfolio

$x_t \in \Delta_n$, asset allocation

$r_t \in \mathbb{R}_+^n$ is price ratio $r_t(i) = \frac{\text{Price}_t(i)}{\text{Price}_t(1)}$

Wealth $W_{t+1} = W_t \cdot r_t^T x_t$

Total wealth $W_T = W_1 \cdot \prod_{t=1}^T r_t^T x_t$
after T days

$$\log \frac{W_T}{W_1} = \sum_{t=1}^T \log(r_t^T x_t) = - \sum_{t=1}^T f_t(x_t)$$

$$f_t(x) = -\log(r_t^T x)$$

$$\text{Regret} = \sum_{t=1}^T f_t(x_t) - \sum_{t=1}^T f_t(u)$$

$u \in \{e_1, e_2, \dots, e_n\}$
 $\uparrow \quad \uparrow \quad \uparrow$
 NVDA AAPL TSLA

$u \in \Delta_n$

$$r_t = \begin{pmatrix} 2 & \frac{1}{2} & 2 & \frac{1}{2} & 2 & \frac{1}{2} & \dots \\ \frac{1}{2} & 2 & \frac{1}{2} & 2 & \frac{1}{2} & 2 & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix}_{t=1, t=2, \dots, t=T}$$

$u = [0.5, 0.5]$

$$\begin{aligned} W_T &= W_1 \cdot \prod_{t=1}^T r_t^T u \\ &= W_1 \cdot \left(\frac{1}{2} \cdot 2 + \frac{1}{2} \cdot \frac{1}{2}\right)^T \\ &\geq W_1 \cdot 1.25^T \end{aligned}$$

OGD: Regret = $C \cdot D \sqrt{T} = O(\sqrt{nT})$

FTRL (entropy Regularizer): $O(\sqrt{T \cdot \log n})$

Can we do better? improve from \sqrt{T}

when f_t is strongly convex: OGD: $\frac{G^2}{m} \log T$

$\hookrightarrow f_t(x) = -\log(r_t^T x)$ strongly convex?

No.

Exponential Concavity:

A convex function f is α -exp-concave if

$e^{-\alpha f(x)}$ is concave.

for all f s.t. f is twice differentiable

f is α -exp-concave at x

$$\Leftrightarrow \nabla^2 f(x) \preceq \alpha \nabla f(x) \nabla f(x)^T$$

Real Strongly Convex

$$\nabla^2 f(x) \succeq m \cdot I$$

Example: $f(x) = -\log(v^T x)$

$$\nabla(-\log(v^T x)) = -\frac{v}{v^T x}$$

$$\nabla^2(-\log(v^T x)) = -\frac{0 - v \cdot v^T}{(v^T x)^2} = \frac{v + v^T}{(v^T x)^2}$$

$$\nabla^2 \succeq \frac{1}{G} \nabla \nabla^T \Rightarrow \text{1-exp-concave.}$$

Example: f is m -strongly convex, in domain K

s.t. $\|\nabla f\|_2 \leq G$

$$\nabla f \nabla f^T \leq G^2 \cdot \mathbf{I} \leq \frac{G^2}{m} \nabla^2 f \quad \text{for all } x \in K$$

\uparrow G -Lipschitz \uparrow m -strongly convex

$$f \text{ is } \nabla^2 f \succeq \frac{m}{G^2} \nabla f \nabla f^T \quad \frac{m}{G^2} \text{-exp-concave.}$$

Useful Property of Exp-Concave Function

(Lemma 4.2 of OCO book Hazan)

f : α -Exp-Concave, domain K , G Lipschitz

$\forall x, y \in K$

$$f(x) \geq f(y) + \nabla f(y)^T (x-y) + \frac{\delta}{2} (x-y)^T \nabla f(y) \nabla f(y)^T (x-y)$$

for $\delta = \frac{1}{2} \min\{\alpha, \frac{1}{4GD}\}$

new for
Exp-concave
functions

Alg. Online Newton Step

Input: $K, T, x_1 \in K, \delta, \epsilon > 0, A_0 = \epsilon \mathbf{I}$

for $t=1, 2, \dots, T$:

1. play x_t

2. incur loss $f_t(x_t)$, receive $\nabla_t = \nabla f_t(x_t)$

3. $A_t = A_{t-1} + \nabla_t \nabla_t^T$

Quasi-Newton
updates

4. Newton Step:

$$y_{t+1} = x_t - \frac{1}{\delta} A_t^{-1} \nabla_t$$

(Projected)

$$x_{t+1} = \Pi_K^{A_t}(y_{t+1}) = \arg \min_{x \in K} \|y_{t+1} - x\|_{A_t}^2$$

Theorem Choose $\gamma = \frac{1}{2} \min\{\alpha, \frac{1}{4GD}\}$
 $\epsilon = \frac{1}{8^2 D^2}$, for all $T \geq 4$
 $\text{Regret}_T \leq 5\left(\frac{1}{\alpha} + GD\right) \frac{n \log T}{T}$

worse \nearrow logarithmic in T

Lemma: $\text{Regret} \leq 4\left(\frac{1}{\alpha} + GD\right) \left(\sum_{t=1}^T \nabla_t^T A_t^{-1} \nabla_t + 1\right)$

Proof: By the 'Useful Property' of α -exp-concave

$$f_t(x^*) \geq f_t(x_t) + \nabla_t^T (x^* - x_t) + \frac{\sigma}{2} (x^* - x_t)^T \nabla_t \nabla_t^T (x^* - x_t)$$

$$f_t(x_t) - f_t(x^*) \leq \nabla_t^T (x_t - x^*) - \frac{\sigma}{2} (x^* - x_t)^T \nabla_t \nabla_t^T (x^* - x_t) \quad (*)$$

Now by the update rule of $x_{t+1} = \Pi_K^{A_t} \left(x_t - \frac{1}{\sigma} A_t^{-1} \nabla_t\right)$

$$\|x_{t+1} - x^*\|_{A_t}^2 \leq \|y_{t+1} - x^*\|_{A_t}^2 = \left\|x_t - \frac{1}{\sigma} A_t^{-1} \nabla_t - x^*\right\|_{A_t}^2$$

non-expansiveness
 i.e. A_t of the
 projection

$$= \|x_t - x^*\|_{A_t}^2 + \frac{1}{\sigma^2} \|A_t^{-1} \nabla_t\|_{A_t}^2 - \frac{2}{\sigma} \underbrace{(x_t - x^*)^T A_t^{-1} \nabla_t}$$

By move it around

$$\frac{\gamma}{2} (x_t - x^*)^T \nabla_t \leq \frac{1}{2\gamma} \|\nabla_t\|_{A_t}^2 + \frac{\gamma}{2} \|x_t - x^*\|_{A_t}^2 - \frac{\gamma}{2} \|x_{t+1} - x^*\|_{A_t}^2$$

multiply $\frac{\gamma}{2}$ on both sides and sum over $t=1, 2, \dots, T$

$$\sum_{t=1}^T (x_t - x^*)^T \nabla_t \leq \frac{1}{2\gamma} \sum_t \|\nabla_t\|_{A_t}^2 + \frac{\gamma}{2} \|x_1 - x^*\|_{A_1}^2 + \sum_{t=2}^T (x_t - x^*)^T (A_t - A_{t-1}) (x_t - x^*) - \frac{\gamma}{2} \|x_{T+1} - x^*\|_{A_T}^2$$

$$\underbrace{A_t = A_{t-1} + \nabla_t \nabla_t^T}_{\text{dropping the negative term}} \leq \frac{1}{2\gamma} \sum_t \|\nabla_t\|_{A_t}^2 + \frac{\gamma}{2} \|x_1 - x^*\|_{A_1}^2 + \sum_{t=2}^T (x_t - x^*)^T \nabla_t \nabla_t^T (x_t - x^*) + 0$$

plug into (*)

$$\sum_t f_t(x_t) - f_t(x^*) \leq \frac{1}{2\gamma} \sum_t \|\nabla_t\|_{A_t}^2 + \frac{\gamma}{2} \|x_1 - x^*\|_{A_1}^2 + \sum_{t=2}^T (x_t - x^*)^T \nabla_t \nabla_t^T (x_t - x^*) \quad \underbrace{\sum_{t=2}^T (x_t - x^*)^T \nabla_t \nabla_t^T (x_t - x^*)}_{\text{blue}}$$

$$= \frac{1}{2\gamma} \sum_t \|\nabla_t\|_{A_t}^2 + \underbrace{\frac{\gamma}{2} (x_1 - x^*)^T A_1 (x_1 - x^*)}_{\text{blue}} - \underbrace{\frac{\gamma}{2} (x_1 - x^*)^T \nabla_1 \nabla_1^T (x_1 - x^*)}_{\text{blue}}$$

$$A_1 = \varepsilon I + \nabla_1 \nabla_1^T \quad = \quad \frac{1}{2\gamma} \sum_t \|\nabla_t\|_{A_t}^2 + \frac{\gamma}{2} (x_1 - x^*)^T (\varepsilon I + \nabla_1 \nabla_1^T) (x_1 - x^*)$$

$$= \frac{1}{2\gamma} \sum_t \|\nabla_t\|_{A_t}^2 + \frac{\gamma}{2} \varepsilon \|x_1 - x^*\|_2^2 \quad \varepsilon = \frac{1}{\gamma^2 D^2}$$

$$\gamma = \frac{1}{2} \min(\alpha, \frac{1}{400})$$

$$= \frac{1}{2\gamma} (\sum_t \|\nabla_t\|_{A_t}^2 + 1)$$

$$\leq \frac{\gamma}{2} \varepsilon \cdot D^2 = \frac{1}{2\gamma}$$

$$\leq 4 \left(\frac{1}{2} + \gamma D \right) \cdot \left(\sum_t \|\nabla_t\|_{A_t}^2 + 1 \right)$$

□

$$\sum_{t=1}^T f_t(x_t) - f_t(x^*) \leq 4\left(\frac{1}{\alpha} + GD\right) \left(\sum_{t=1}^T \underbrace{\|\nabla_t\|_{A_t^{-1}}^2}_{\text{Remains to bound this term!}} + 1 \right)$$

Remains to bound this term!

Proof of the main theorem:

$$\|\nabla_t\|_{A_t^{-1}}^2 = \nabla_t^T A_t^{-1} \nabla_t = \text{tr}(\nabla_t^T A_t^{-1} \nabla_t) = \text{tr}(\underbrace{A_t^{-1}}_{\mathbb{R}^{n \times n}} \underbrace{\nabla_t \nabla_t^T}_{\mathbb{R}^{n \times n}}) =$$

$$A_t^{-1} \bullet (\nabla_t \nabla_t^T)$$

$$= \text{tr}(A_t^{-1} (A_t - A_{t-1})) = \text{tr}(I - A_t^{-1} A_{t-1}) = \sum_{i=1}^n (1 - \lambda_i(A_t^{-1} A_{t-1}))$$

$$\underbrace{\varepsilon \leq \log\left(\frac{1}{1-\varepsilon}\right)}_{\text{take}} \leq \sum_{i=1}^n \log(\lambda_i^{-1}(A_t^{-1} A_{t-1})) = \log\left(\prod_{i=1}^n \lambda_i^{-1}(A_t^{-1} A_{t-1})\right) = \log\left|(A_t^{-1} A_{t-1})^{-1}\right|$$

$$\varepsilon = 1 - \lambda_i(A_t^{-1} A_{t-1})$$

$$= \log\left(\frac{|A_t|}{|A_{t-1}|}\right) = \log|A_t| - \log|A_{t-1}|$$

add up $t=1, 2, \dots, T$

$$\sum_{t=1}^T \|\nabla_t\|_{A_t^{-1}}^2 = \log|A_T| - \log|A_0|$$

$$A_0 = \varepsilon I, \quad |A_0| = \varepsilon^n$$

$$A_T = \sum_{t=1}^T \nabla_t \nabla_t^T + \varepsilon I \preceq (TG^2 + \varepsilon) I$$

$$|A_T| \leq (TG^2 + \varepsilon)^n$$

$$\gamma = \frac{1}{2} \min\left(\alpha, \frac{1}{4GD}\right) \leq n \cdot \log\left(\frac{1}{\varepsilon} (TG^2 + 1)\right)$$

$$\varepsilon = \frac{1}{\gamma^2 D^2} = n \cdot \log\left(\gamma^2 T G^2 D^2 + 1\right) \leq n \cdot \log\left(\frac{1}{\alpha^2} T \frac{G^2 D^2}{\alpha^2} + 1\right) \leq n \cdot \log T \text{ for } T \geq 4. \quad \square$$

Complexity (computational)

$A_t^{-1} \nabla_t$ requires $O(n^3)$

space
 $O(n^2)$

Scherman-Morrison-Woodbury (Matrix inversion lemma)

$$(A + xx^T)^{-1} = A^{-1} - \frac{A^{-1}xx^T A^{-1}}{1 + x^T A^{-1}x}$$

Newton step from previous iteration

$O(n^2)$

Projection $\min_{x \in K} \|y_{t+1} - x\|_{A_t}^2$ $O(n^2 \log \frac{1}{\epsilon})$
 $\hookrightarrow O(n^2)$

$OAD = O(n)$ space and time