

Online and Stochastic Gradient Descent, Multi-Armed Bandits

Dimitris Fotakis

SCHOOL OF ELECTRICAL AND COMPUTER ENGINEERING
NATIONAL TECHNICAL UNIVERSITY OF ATHENS, GREECE

Online Gradient Descent

Online Gradient Descent

Input: convex set S , $w_1 \in S$, time horizon T , step size η

For each $t = 1, \dots, T$ do:

- Play w_t , get cost function $f_t : S \rightarrow \mathbb{R}$, incur cost $f_t(w_t)$
- Update $y_{t+1} = w_t - \eta \nabla f_t(w_t)$
- Project $w_{t+1} = \arg \min_{w \in S} \|w - y_{t+1}\|$

Online Gradient Descent

Online Gradient Descent

Input: convex set S , $w_1 \in S$, time horizon T , step size η

For each $t = 1, \dots, T$ do:

- Play w_t , get cost function $f_t : S \rightarrow \mathbb{R}$, incur cost $f_t(w_t)$
- Update $y_{t+1} = w_t - \eta \nabla f_t(w_t)$
- Project $w_{t+1} = \arg \min_{w \in S} \|w - y_{t+1}\|$

Ignore project step for the analysis: $w_{t+1} = w_t - \eta \nabla f_t(w_t)$

Also we let $v_t = \nabla f_t(w_t)$ for brevity.

Online Gradient Descent

Online Gradient Descent

Input: convex set S , $w_1 \in S$, time horizon T , step size η

For each $t = 1, \dots, T$ do:

- Play w_t , get cost function $f_t : S \rightarrow \mathbb{R}$, incur cost $f_t(w_t)$
- Update $y_{t+1} = w_t - \eta \nabla f_t(w_t)$
- Project $w_{t+1} = \arg \min_{w \in S} \|w - y_{t+1}\|$

Ignore project step for the analysis: $w_{t+1} = w_t - \eta \nabla f_t(w_t)$

Also we let $v_t = \nabla f_t(w_t)$ for brevity.

Using convexity $f_t(w^*) \geq f_t(w_t) + \nabla f_t(w_t)(w^* - w_t)$, we get that:

$$\sum_{t=1}^T (f_t(w_t) - f_t(w^*)) \leq \sum_{t=1}^T v_t(w_t - w^*) \quad (1)$$

Online Gradient Descent

From the update rule, we get that:

$$\begin{aligned}\|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2 &= \|\mathbf{w}_t - \eta \mathbf{v}_t - \mathbf{w}^*\|^2 \\ &= \|\mathbf{w}_t - \mathbf{w}^*\|^2 - 2\eta \mathbf{v}_t(\mathbf{w}_t - \mathbf{w}^*) + \eta^2 \|\mathbf{v}_t\|^2 \Rightarrow \\ \mathbf{v}_t(\mathbf{w}_t - \mathbf{w}^*) &= \frac{1}{2\eta} (\|\mathbf{w}_t - \mathbf{w}^*\|^2 - \|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2) + \frac{\eta}{2} \|\mathbf{v}_t\|^2\end{aligned}\quad (2)$$

From the update rule, we get that:

$$\begin{aligned}\|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2 &= \|\mathbf{w}_t - \eta \mathbf{v}_t - \mathbf{w}^*\|^2 \\ &= \|\mathbf{w}_t - \mathbf{w}^*\|^2 - 2\eta \mathbf{v}_t(\mathbf{w}_t - \mathbf{w}^*) + \eta^2 \|\mathbf{v}_t\|^2 \Rightarrow \\ \mathbf{v}_t(\mathbf{w}_t - \mathbf{w}^*) &= \frac{1}{2\eta} (\|\mathbf{w}_t - \mathbf{w}^*\|^2 - \|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2) + \frac{\eta}{2} \|\mathbf{v}_t\|^2\end{aligned}\quad (2)$$

Substituting (2) in (1) results in a telescopic sum. So, we get that:

$$\begin{aligned}\text{Regret}_{\text{OGD}} &= \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{w}^*)) \leq \frac{\|\mathbf{w}_1 - \mathbf{w}^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\mathbf{v}_t\|^2 \\ &\leq \frac{B^2}{2\eta} + \frac{\eta T G^2}{2} \stackrel{\eta = \frac{B}{G\sqrt{T}}}{=} B G \sqrt{T}\end{aligned}$$

From the update rule, we get that:

$$\begin{aligned}\|w_{t+1} - w^*\|^2 &= \|w_t - \eta v_t - w^*\|^2 \\ &= \|w_t - w^*\|^2 - 2\eta v_t(w_t - w^*) + \eta^2 \|v_t\|^2 \Rightarrow \\ v_t(w_t - w^*) &= \frac{1}{2\eta} (\|w_t - w^*\|^2 - \|w_{t+1} - w^*\|^2) + \frac{\eta}{2} \|v_t\|^2\end{aligned}\quad (2)$$

Substituting (2) in (1) results in a telescopic sum. So, we get that:

$$\begin{aligned}\text{Regret}_{OGD} &= \sum_{t=1}^T (f_t(w_t) - f_t(w^*)) \leq \frac{\|w_1 - w^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2 \\ &\leq \frac{B^2}{2\eta} + \frac{\eta TG^2}{2} \stackrel{\eta = \frac{B}{G\sqrt{T}}}{=} BG\sqrt{T}\end{aligned}$$

Similar regret with step $\eta_t = \frac{B}{G\sqrt{t}}$, because $\sum_{t=1}^T \frac{1}{\sqrt{t}} \leq 2\sqrt{T}$.

α -Strongly Convex Functions

Using α -strong convexity

$$f_t(\mathbf{w}^*) \geq f_t(\mathbf{w}_t) + \nabla f_t(\mathbf{w}_t)(\mathbf{w}^* - \mathbf{w}_t) + \frac{\alpha}{2} \|\mathbf{w}^* - \mathbf{w}_t\|^2,$$

we get that:

$$2 \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{w}^*)) \leq \sum_{t=1}^T \left(2v_t(\mathbf{w}_t - \mathbf{w}^*) - \alpha \|\mathbf{w}_t - \mathbf{w}^*\|^2 \right) \quad (3)$$

α -Strongly Convex Functions

Using α -strong convexity

$$f_t(\mathbf{w}^*) \geq f_t(\mathbf{w}_t) + \nabla f_t(\mathbf{w}_t)(\mathbf{w}^* - \mathbf{w}_t) + \frac{\alpha}{2} \|\mathbf{w}^* - \mathbf{w}_t\|^2,$$

we get that:

$$2 \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{w}^*)) \leq \sum_{t=1}^T \left(2\mathbf{v}_t(\mathbf{w}_t - \mathbf{w}^*) - \alpha \|\mathbf{w}_t - \mathbf{w}^*\|^2 \right) \quad (3)$$

Substituting (2) in (3), we get that:

$$\begin{aligned} 2\text{Regret}_{\text{OGD}} &= 2 \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{w}^*)) \\ &\leq \sum_{t=1}^T \|\mathbf{w}_t - \mathbf{w}^*\|^2 \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} - \alpha \right) + \sum_{t=1}^T \eta_t \|\mathbf{v}_t\|^2 \\ &\stackrel{\eta_t=1/(\alpha t)}{\leq} 0 + \frac{G^2(1 + \ln T)}{\alpha} \end{aligned}$$

Learning Using Stochastic Gradient Descent

- In learning problems, we want to solve $\min_{\mathbf{w} \in H} L_{\mathcal{D}}(\mathbf{w})$, where $L_{\mathcal{D}}(H) = \mathbb{E}_{z \sim \mathcal{D}}[\ell(\mathbf{w}, z)]$ and $\ell(\mathbf{w}, z)$ is the loss on z under $\mathbf{w} \in H$.
- Instead of Empirical Risk Minimization $\min_{\mathbf{w} \in H} L_S(\mathbf{w})$, we try to minimize $L_{\mathcal{D}}(\mathbf{w})$ directly using SGD and the fact that:

$$\nabla L_{\mathcal{D}}(\mathbf{w}) = \nabla \mathbb{E}_{z \sim \mathcal{D}}[\ell(\mathbf{w}, z)] = \mathbb{E}_{z \sim \mathcal{D}}[\nabla \ell(\mathbf{w}, z)]$$

Learning Using Stochastic Gradient Descent

- In learning problems, we want to solve $\min_{\mathbf{w} \in H} L_{\mathcal{D}}(\mathbf{w})$, where $L_{\mathcal{D}}(H) = \mathbb{E}_{z \sim \mathcal{D}}[\ell(\mathbf{w}, z)]$ and $\ell(\mathbf{w}, z)$ is the loss on z under $\mathbf{w} \in H$.
- Instead of Empirical Risk Minimization $\min_{\mathbf{w} \in H} L_S(\mathbf{w})$, we try to minimize $L_{\mathcal{D}}(\mathbf{w})$ directly using SGD and the fact that:

$$\nabla L_{\mathcal{D}}(\mathbf{w}) = \nabla \mathbb{E}_{z \sim \mathcal{D}}[\ell(\mathbf{w}, z)] = \mathbb{E}_{z \sim \mathcal{D}}[\nabla \ell(\mathbf{w}, z)]$$

- We cannot calculate $\nabla L_{\mathcal{D}}(\mathbf{w})$, because we do not know \mathcal{D} .
- But expected value of $\nabla \ell(\mathbf{w}, z)$, over $z \sim \mathcal{D}$, is $\nabla L_{\mathcal{D}}(\mathbf{w})$
- We use $\nabla \ell(\mathbf{w}, z)$, for sample $z \sim \mathcal{D}$, which is an **unbiased estimator** of the gradient.

Stochastic Gradient Descent

Stochastic Gradient Descent

Input: convex set H , $\mathbf{w}_1 \in H$, time horizon T , step size η

For each $t = 1, \dots, T$ do:

- Random sample $z_t \sim \mathcal{D}$, $\mathbf{v}_t = \nabla \ell(\mathbf{w}_t, z_t)$ and $f_t(\mathbf{w}_t) = \mathbf{v}_t \mathbf{w}_t$
- Update $\mathbf{y}_{t+1} = \mathbf{w}_t - \eta \mathbf{v}_t$
- Project $\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in H} \|\mathbf{w} - \mathbf{y}_{t+1}\|$

Return $\bar{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$

Stochastic Gradient Descent

Stochastic Gradient Descent

Input: convex set H , $\mathbf{w}_1 \in H$, time horizon T , step size η

For each $t = 1, \dots, T$ do:

- Random sample $z_t \sim \mathcal{D}$, $\mathbf{v}_t = \nabla \ell(\mathbf{w}_t, z_t)$ and $f_t(\mathbf{w}_t) = \mathbf{v}_t \mathbf{w}_t$
- Update $\mathbf{y}_{t+1} = \mathbf{w}_t - \eta \mathbf{v}_t$
- Project $\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in H} \|\mathbf{w} - \mathbf{y}_{t+1}\|$

Return $\bar{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$

Analysis based on OGD, with $f_t(\mathbf{w}) = \mathbf{v}_t \mathbf{w}$.

- SGD tries to minimize $f(\mathbf{w}) = L_{\mathcal{D}}(\mathbf{w})$
- Note that $\nabla f(\mathbf{w}_t) = \mathbb{E}_{z \sim \mathcal{D}} [\nabla \ell(\mathbf{w}_t, z)] = \mathbb{E}_{z \sim \mathcal{D}} [\mathbf{v}_t] = \nabla f_t(\mathbf{w}_t)$

Stochastic Gradient Descent: Analysis

Analysis based on OGD, with $f_t(\mathbf{w}) = \mathbf{v}_t \mathbf{w}$

- SGD tries to minimize $f(\mathbf{w}) = L_{\mathcal{D}}(\mathbf{w})$
- Note that $\nabla f(\mathbf{w}_t) = \mathbb{E}_{\mathbf{z} \sim \mathcal{D}}[\nabla \ell(\mathbf{w}_t, \mathbf{z})] = \mathbb{E}_{\mathbf{z} \sim \mathcal{D}}[\mathbf{v}_t] = \nabla f_t(\mathbf{w}_t)$

$$\begin{aligned} \mathbb{E}_{z_1, \dots, z_T} [f(\bar{\mathbf{w}})] - f(\mathbf{w}^*) &\leq \mathbb{E} \left[\sum_{t=1}^T \frac{f(\mathbf{w}_t) - f(\mathbf{w}^*)}{T} \right] \\ &\leq \frac{1}{T} \mathbb{E} \left[\sum_{t=1}^T \nabla f(\mathbf{w}_t) (\mathbf{w}_t - \mathbf{w}^*) \right] \\ &= \frac{1}{T} \mathbb{E} \left[\sum_{t=1}^T \mathbf{v}_t (\mathbf{w}_t - \mathbf{w}^*) \right] \\ &= \frac{1}{T} \mathbb{E} \left[\sum_{t=1}^T f_t(\mathbf{w}_t) - f_t(\mathbf{w}^*) \right] \\ &\leq \frac{\text{Regret}_{\text{OGD}}(T)}{T} \leq \frac{BG}{\sqrt{T}} \end{aligned}$$