

# Online Convex Optimization, Online and Stochastic Gradient Descent

**Dimitris Fotakis**

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

# Online Convex Optimization

General framework: convex set  $S \subseteq \mathbb{R}^d$

On each day  $t = 1, \dots, T$ :

- ➊ Learner picks vector  $p_t \in S$
- ➋ Adversary picks **convex loss** function  $f_t : S \rightarrow \mathbb{R}$ ,  
with  $f_t$  differentiable and  $L$ -Lipschitz wrt some norm  $\|\cdot\|$ ,  
i.e.,  $|f_t(p) - f_t(p')| \leq L \cdot \|p - p'\|$
- ➌ Learner **learns  $f_t$**  and incurs loss  $f_t(p_t)$

# Online Convex Optimization

General framework: convex set  $S \subseteq \mathbb{R}^d$

On each day  $t = 1, \dots, T$ :

- ➊ Learner picks vector  $p_t \in S$
- ➋ Adversary picks **convex loss** function  $f_t : S \rightarrow \mathbb{R}$ ,  
with  $f_t$  differentiable and  $L$ -Lipschitz wrt some norm  $\|\cdot\|$ ,  
i.e.,  $|f_t(p) - f_t(p')| \leq L \cdot \|p - p'\|$
- ➌ Learner **learns  $f_t$**  and incurs loss  $f_t(p_t)$

Goal is to minimize **regret**:

$$\text{Regret}(T) = \sup_{f_1, \dots, f_T} \left( \sum_{t=1}^T f_t(p_t) - \min_{p \in S} \sum_{t=1}^T f_t(p) \right)$$

(Online) algorithm is **no-regret** if  $\text{Regret}(T)/T \rightarrow 0$  at  $T \rightarrow \infty$

# Online Convex Optimization

General framework: convex set  $S \subseteq \mathbb{R}^d$

On each day  $t = 1, \dots, T$ :

- ➊ Learner picks vector  $p_t \in S$
- ➋ Adversary picks **convex loss** function  $f_t : S \rightarrow \mathbb{R}$ ,  
with  $f_t$  differentiable and  $L$ -Lipschitz wrt some norm  $\|\cdot\|$ ,  
i.e.,  $|f_t(p) - f_t(p')| \leq L \cdot \|p - p'\|$
- ➌ Learner **learns  $f_t$**  and incurs loss  $f_t(p_t)$

Goal is to minimize **regret**:

$$\text{Regret}(T) = \sup_{f_1, \dots, f_T} \left( \sum_{t=1}^T f_t(p_t) - \min_{p \in S} \sum_{t=1}^T f_t(p) \right)$$

(Online) algorithm is **no-regret** if  $\text{Regret}(T)/T \rightarrow 0$  at  $T \rightarrow \infty$

- **Experts**:  $S = \Delta_d$ ,  $f_t(p) = \langle p, \ell_t \rangle$  (linear expected loss)

# Online Convex Optimization

General framework: convex set  $S \subseteq \mathbb{R}^d$

On each day  $t = 1, \dots, T$ :

- ➊ Learner picks vector  $p_t \in S$
- ➋ Adversary picks **convex loss** function  $f_t : S \rightarrow \mathbb{R}$ ,  
with  $f_t$  differentiable and  $L$ -Lipschitz wrt some norm  $\|\cdot\|$ ,  
i.e.,  $|f_t(p) - f_t(p')| \leq L \cdot \|p - p'\|$
- ➌ Learner **learns  $f_t$**  and incurs loss  $f_t(p_t)$

Goal is to minimize **regret**:

$$\text{Regret}(T) = \sup_{f_1, \dots, f_T} \left( \sum_{t=1}^T f_t(p_t) - \min_{p \in S} \sum_{t=1}^T f_t(p) \right)$$

(Online) algorithm is **no-regret** if  $\text{Regret}(T)/T \rightarrow 0$  at  $T \rightarrow \infty$

- **Experts**:  $S = \Delta_d, f_t(p) = \langle p, \ell_t \rangle$  (linear expected loss)
- **Online Quadratic Optimization**: learner  $p_t \in S$ ,  
adversary  $z_t \in S, f_t(p) = \|p - z_t\|_2^2$ .

# Online Convex Optimization

General framework: convex set  $S \subseteq \mathbb{R}^d$

On each day  $t = 1, \dots, T$ :

- ➊ Learner picks vector  $p_t \in S$
- ➋ Adversary picks **convex loss** function  $f_t : S \rightarrow \mathbb{R}$ ,  
with  $f_t$  differentiable and  $L$ -Lipschitz wrt some norm  $\|\cdot\|$ ,  
i.e.,  $|f_t(p) - f_t(p')| \leq L \cdot \|p - p'\|$
- ➌ Learner **learns  $f_t$**  and incurs loss  $f_t(p_t)$

Goal is to minimize **regret**:

$$\text{Regret}(T) = \sup_{f_1, \dots, f_T} \left( \sum_{t=1}^T f_t(p_t) - \min_{p \in S} \sum_{t=1}^T f_t(p) \right)$$

(Online) algorithm is **no-regret** if  $\text{Regret}(T)/T \rightarrow 0$  at  $T \rightarrow \infty$

- **Experts**:  $S = \Delta_d, f_t(p) = \langle p, \ell_t \rangle$  (linear expected loss)
- **Online Quadratic Optimization**: learner  $p_t \in S$ ,  
adversary  $z_t \in S, f_t(p) = \|p - z_t\|_2^2$ .
- **Online Least Squares Linear Regression**: learner  $p_t \in \mathbb{R}^d$ ,  
 $\|p_t\| \leq B$ , adversary  $(x_t, y_t) \in \mathbb{R}^d \times \mathbb{R}, f_t(p) = (\langle p, x_t \rangle - y_t)^2$

# Follow / Be the Regularized Leader

$$F_t(p) = \sum_{\tau=1}^t f_\tau(p) \text{ and } \tilde{F}_t(p) = \sum_{\tau=1}^t f_\tau(p) + R(p)/\eta$$

**FTRL**:  $\tilde{p}_t = \arg \min_{p \in S} \tilde{F}_{t-1}(p)$

**BTRL**:  $\tilde{p}_t^* = \arg \min_{p \in S} \tilde{F}_t(p)$

# Follow / Be the Regularized Leader

$$F_t(p) = \sum_{\tau=1}^t f_\tau(p) \text{ and } \tilde{F}_t(p) = \sum_{\tau=1}^t f_\tau(p) + R(p)/\eta$$

**FTRL**:  $\tilde{p}_t = \arg \min_{p \in S} \tilde{F}_{t-1}(p)$

**BTRL**:  $\tilde{p}_t^* = \arg \min_{p \in S} \tilde{F}_t(p)$

$1/\eta$ -strongly convex function  $R : S \rightarrow \mathbb{R}$  wrt norm  $\|\cdot\|$ , if  $\forall x, y \in S$ :

$$R(x) \geq R(y) + \langle \nabla R(y), x - y \rangle + \frac{1}{2\eta} \|x - y\|^2$$

# Follow / Be the Regularized Leader

$$F_t(p) = \sum_{\tau=1}^t f_\tau(p) \text{ and } \tilde{F}_t(p) = \sum_{\tau=1}^t f_\tau(p) + R(p)/\eta$$

**FTRL**:  $\tilde{p}_t = \arg \min_{p \in S} \tilde{F}_{t-1}(p)$

**BTRL**:  $\tilde{p}_t^* = \arg \min_{p \in S} \tilde{F}_t(p)$

**1/ $\eta$ -strongly convex** function  $R : S \rightarrow \mathbb{R}$  wrt norm  $\|\cdot\|$ , if  $\forall x, y \in S$ :

$$R(x) \geq R(y) + \langle \nabla R(y), x - y \rangle + \frac{1}{2\eta} \|x - y\|^2$$

Functions  $f, g : S \rightarrow \mathbb{R}$  be **1/ $\eta$ -strongly convex** wrt some norm  $\|\cdot\|$  and  $h(x) = g(x) - f(x)$  be **L-Lipschitz** wrt  $\|\cdot\|$ .

Then,  $\|x_f^* - x_g^*\| \leq \eta \cdot L$ , with  $x_f^*, x_g^*$  **minimizers** of  $f, g$ .

# Regret of FTRL Against BTRL

$$\begin{aligned}\text{Regret}_{FTRL}(T) &\leq \text{Regret}_{BTRL}(T) + L \cdot \sum_{t=1}^T \|\tilde{p}_t - \tilde{p}_{t+1}\| \\ &\leq \text{Regret}_{BTRL}(T) + \eta \cdot L^2 \cdot T\end{aligned}$$

# Regret of FTRL Against BTRL

$$\begin{aligned}\text{Regret}_{FTRL}(T) &\leq \text{Regret}_{BTRL}(T) + L \cdot \sum_{t=1}^T \|\tilde{p}_t - \tilde{p}_{t+1}\| \\ &\leq \text{Regret}_{BTRL}(T) + \eta \cdot L^2 \cdot T\end{aligned}$$

**Proof:** Second inequality from strong convexity, because  $\tilde{p}_t, \tilde{p}_{t+1}$  are minimizers of  $1/\eta$ -strong convex functions  $\tilde{F}_{t-1}(p)$  and  $\tilde{F}_t(p)$  with difference  $f_t(p)$  which is  $L$ -Lipschitz.

# Regret of FTRL Against BTRL

$$\begin{aligned}\text{Regret}_{FTRL}(T) &\leq \text{Regret}_{BTRL}(T) + L \cdot \sum_{t=1}^T \|\tilde{p}_t - \tilde{p}_{t+1}\| \\ &\leq \text{Regret}_{BTRL}(T) + \eta \cdot L^2 \cdot T\end{aligned}$$

**Proof:** Second inequality from strong convexity, because  $\tilde{p}_t, \tilde{p}_{t+1}$  are minimizers of  $1/\eta$ -strong convex functions  $\tilde{F}_{t-1}(p)$  and  $\tilde{F}_t(p)$  with difference  $f_t(p)$  which is  $L$ -Lipschitz.

For the first inequality, we observe that:

$$\begin{aligned}\text{Regret}_{FTRL}(T) - \text{Regret}_{BTRL}(T) &= \sum_{t=1}^T (f_t(\tilde{p}_t) - f_t(\tilde{p}_t^*)) \\ &\leq L \sum_{t=1}^T \|\tilde{p}_t - \tilde{p}_t^*\| \\ &= L \sum_{t=1}^T \|\tilde{p}_t - \tilde{p}_{t+1}\|\end{aligned}$$

# Regret of Be the Regularized Leader

$$\text{Regret}_{BTRL}(T) \leq \frac{1}{\eta} \left( \max_{p \in S} R(p) - \min_{p \in S} R(p) \right)$$

# Regret of Be the Regularized Leader

$$\text{Regret}_{BTRL}(T) \leq \frac{1}{\eta} \left( \max_{p \in S} R(p) - \min_{p \in S} R(p) \right)$$

**Proof:**

- Let  $f_0(p) = R(p)/\eta$  and  $\tilde{p}_0^* = \arg \min_{p \in S} R(p)/\eta$ .
- Using induction on  $t$ , we show that for all  $t \geq 0$ ,

$$\sum_{\tau=0}^t f_\tau(\tilde{p}_\tau^*) \leq \min_{p \in S} \sum_{\tau=0}^t f_\tau(p) \quad (\text{notice fake action } \tilde{p}_0^* \text{ at } \tau = 0)$$

# Regret of Be the Regularized Leader

$$\text{Regret}_{BTRL}(T) \leq \frac{1}{\eta} \left( \max_{p \in S} R(p) - \min_{p \in S} R(p) \right)$$

**Proof:**

- Let  $f_0(p) = R(p)/\eta$  and  $\tilde{p}_0^* = \arg \min_{p \in S} R(p)/\eta$ .
- Using induction on  $t$ , we show that for all  $t \geq 0$ ,

$$\sum_{\tau=0}^t f_\tau(\tilde{p}_\tau^*) \leq \min_{p \in S} \sum_{\tau=0}^t f_\tau(p) \quad (\text{notice fake action } \tilde{p}_0^* \text{ at } \tau = 0)$$

- Then, using the claim above,

$$\sum_{t=0}^T f_t(\tilde{p}_t^*) \leq \min_{p \in S} \sum_{t=0}^T f_t(p) \leq \max_{p \in S} f_0(p) + \min_{p \in S} \sum_{t=1}^T f_t(p)$$

# Regret of Be the Regularized Leader

$$\text{Regret}_{BTRL}(T) \leq \frac{1}{\eta} \left( \max_{p \in S} R(p) - \min_{p \in S} R(p) \right)$$

**Proof:**

- Let  $f_0(p) = R(p)/\eta$  and  $\tilde{p}_0^* = \arg \min_{p \in S} R(p)/\eta$ .
- Using induction on  $t$ , we show that for all  $t \geq 0$ ,

$$\sum_{\tau=0}^t f_\tau(\tilde{p}_\tau^*) \leq \min_{p \in S} \sum_{\tau=0}^t f_\tau(p) \quad (\text{notice fake action } \tilde{p}_0^* \text{ at } \tau = 0)$$

- Then, using the claim above,

$$\sum_{t=0}^T f_t(\tilde{p}_t^*) \leq \min_{p \in S} \sum_{t=0}^T f_t(p) \leq \max_{p \in S} f_0(p) + \min_{p \in S} \sum_{t=1}^T f_t(p)$$

- Hence, by rearranging:

$$\sum_{t=1}^T f_t(\tilde{p}_t^*) - \min_{p \in S} \sum_{t=1}^T f_t(p) \leq \max_{p \in S} R(p)/\eta - \min_{p \in S} R(p)/\eta$$

# Regret of Follow the Regularized Leader

**Theorem :**

$$\text{Regret}_{FTRL}(T) \leq \eta \cdot L^2 \cdot T + \frac{(\max_{p \in S} R(p) - \min_{p \in S} R(p))}{\eta}$$

# Regret of Follow the Regularized Leader

**Theorem :**

$$\text{Regret}_{FTRL}(T) \leq \eta \cdot L^2 \cdot T + \frac{(\max_{p \in S} R(p) - \min_{p \in S} R(p))}{\eta}$$

Let  $R^* = \max_{p \in S} R(p) - \min_{p \in S} R(p)$ .

Setting  $\eta = \sqrt{R^*/T}$  yields  $\text{Regret}_{FTRL}(T) \leq (L^2 + 1)\sqrt{R^*T}$

# Regret of Follow the Regularized Leader

**Theorem :**

$$\text{Regret}_{FTRL}(T) \leq \eta \cdot L^2 \cdot T + \frac{(\max_{p \in S} R(p) - \min_{p \in S} R(p))}{\eta}$$

Let  $R^* = \max_{p \in S} R(p) - \min_{p \in S} R(p)$ .

Setting  $\eta = \sqrt{R^*/T}$  yields  $\text{Regret}_{FTRL}(T) \leq (L^2 + 1)\sqrt{R^*T}$

**Multiplicative Weight Updates :**

- Negative entropy  $E^-(p) = \sum_{i=1}^d p_i \ln(p_i)$  is 1-strongly convex wrt  $L_1$  norm.
- Using  $E^-(p)$  as regularizer, results in the following update rule for linear losses  $f_t(p) = \langle p, \ell_t \rangle$ :

$$p_{t+1}(i) = p_t(i) \cdot e^{-\eta \ell_t(i)} \approx p_t(i)(1 - \eta \ell_t(i))$$

# Regret of Follow the Regularized Leader

**Theorem :**

$$\text{Regret}_{FTRL}(T) \leq \eta \cdot L^2 \cdot T + \frac{(\max_{p \in S} R(p) - \min_{p \in S} R(p))}{\eta}$$

Let  $R^* = \max_{p \in S} R(p) - \min_{p \in S} R(p)$ .

Setting  $\eta = \sqrt{R^*/T}$  yields  $\text{Regret}_{FTRL}(T) \leq (L^2 + 1)\sqrt{R^*T}$

**Multiplicative Weight Updates :**

- Negative entropy  $E^-(p) = \sum_{i=1}^d p_i \ln(p_i)$  is 1-strongly convex wrt  $L_1$  norm.
- Using  $E^-(p)$  as regularizer, results in the following update rule for linear losses  $f_t(p) = \langle p, \ell_t \rangle$ :

$$p_{t+1}(i) = p_t(i) \cdot e^{-\eta \ell_t(i)} \approx p_t(i)(1 - \eta \ell_t(i))$$

- If  $\ell_t \in [0, 1]^d$ , setting  $\eta = \sqrt{\ln(d)/T}$ , yields regret  $2\sqrt{T \ln(d)}$

## Online Projected Gradient Descent

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

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

- Play  $\mathbf{w}_t$ , get **convex cost** function  $f_t : S \rightarrow \mathbb{R}$ , incur cost  $f_t(\mathbf{w}_t)$
- Update  $\mathbf{y}_{t+1} = \mathbf{w}_t - \eta \nabla f_t(\mathbf{w}_t)$
- Project  $\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in S} \|\mathbf{w} - \mathbf{y}_{t+1}\|$

## Online Projected Gradient Descent

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

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

- Play  $\mathbf{w}_t$ , get **convex cost** function  $f_t : S \rightarrow \mathbb{R}$ , incur cost  $f_t(\mathbf{w}_t)$
- Update  $\mathbf{y}_{t+1} = \mathbf{w}_t - \eta \nabla f_t(\mathbf{w}_t)$
- Project  $\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in S} \|\mathbf{w} - \mathbf{y}_{t+1}\|$

**Ignore project** step for the analysis:  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla f_t(\mathbf{w}_t)$

Also we let  $\mathbf{v}_t = \nabla f_t(\mathbf{w}_t)$  for brevity.

## Online Projected Gradient Descent

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

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

- Play  $\mathbf{w}_t$ , get **convex cost** function  $f_t : S \rightarrow \mathbb{R}$ , incur cost  $f_t(\mathbf{w}_t)$
- Update  $\mathbf{y}_{t+1} = \mathbf{w}_t - \eta \nabla f_t(\mathbf{w}_t)$
- Project  $\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in S} \|\mathbf{w} - \mathbf{y}_{t+1}\|$

**Ignore project** step for the analysis:  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla f_t(\mathbf{w}_t)$

Also we let  $\mathbf{v}_t = \nabla f_t(\mathbf{w}_t)$  for brevity.

Using convexity  $f_t(\mathbf{w}^*) \geq f(\mathbf{w}_t) + \nabla f_t(\mathbf{w}_t)(\mathbf{w}^* - \mathbf{w}_t)$ , we get that:

$$\text{Regret}_{OGD} = \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{w}^*)) \leq \sum_{t=1}^T \mathbf{v}_t (\mathbf{w}_t - \mathbf{w}^*) \quad (1)$$

# Online Gradient Descent

From the update rule, we get that:

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

# Online Gradient Descent

From the update rule, we get that:

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

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

# Online Gradient Descent

From the update rule, we get that:

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

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

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

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

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

# $\alpha$ -Strongly Convex Functions

Using  $\alpha$ -strong convexity

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

we get that:

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

# $\alpha$ -Strongly Convex Functions

Using  $\alpha$ -strong convexity

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

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:

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

# Learning Using Stochastic Gradient Descent

- In learning problems, we want to solve  $\min_{\mathbf{w} \in H} L_{\mathcal{D}}(\mathbf{w})$ , where  $L_{\mathcal{D}}(\mathbf{w}) = \mathbb{E}\text{Exp}_{z \sim \mathcal{D}}[\ell(\mathbf{w}, z)]$  and  $\ell(\mathbf{w}, z)$  is the loss on  $z$  under  $\mathbf{w} \in H$ .

# Learning Using Stochastic Gradient Descent

- In learning problems, we want to solve  $\min_{\mathbf{w} \in H} L_{\mathcal{D}}(\mathbf{w})$ , where  $L_{\mathcal{D}}(\mathbf{w}) = \mathbb{E}\text{Exp}_{z \sim \mathcal{D}}[\ell(\mathbf{w}, z)]$  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:

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

# Learning Using Stochastic Gradient Descent

- In learning problems, we want to solve  $\min_{\mathbf{w} \in H} L_{\mathcal{D}}(\mathbf{w})$ , where  $L_{\mathcal{D}}(\mathbf{w}) = \mathbb{E}\text{Exp}_{z \sim \mathcal{D}}[\ell(\mathbf{w}, z)]$  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:

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

- 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  $\eta$

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  $f_t(\mathbf{w}) = \mathbf{v}_t \mathbf{w}$
- 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  $\eta$

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  $f_t(\mathbf{w}) = \mathbf{v}_t \mathbf{w}$
- 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  $\nabla f(\mathbf{w}_t) = \mathbb{E}_{z \sim \mathcal{D}} [\nabla \ell(\mathbf{w}_t, z)] = \mathbb{E}_{z \sim \mathcal{D}} [\mathbf{v}_t] = \nabla f_t(\mathbf{w}_t)$

$$\mathbb{E}_{z_1, \dots, z_T} [f(\bar{\mathbf{w}})] - f(\mathbf{w}^*) \leq \mathbb{E} \left[ \sum_{t=1}^T \frac{f(\mathbf{w}_t) - f(\mathbf{w}^*)}{T} \right]$$

# 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  $\nabla f(\mathbf{w}_t) = \mathbb{E} \text{Exp}_{z \sim \mathcal{D}} [\nabla \ell(\mathbf{w}_t, z)] = \mathbb{E} \text{Exp}_{z \sim \mathcal{D}} [\mathbf{v}_t] = \nabla f_t(\mathbf{w}_t)$

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

# 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  $\nabla f(\mathbf{w}_t) = \mathbb{E}_{z \sim \mathcal{D}} [\nabla \ell(\mathbf{w}_t, z)] = \mathbb{E}_{z \sim \mathcal{D}} [\mathbf{v}_t] = \nabla f_t(\mathbf{w}_t)$

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

# 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  $\nabla f(\mathbf{w}_t) = \mathbb{E} \text{Exp}_{z \sim \mathcal{D}} [\nabla \ell(\mathbf{w}_t, z)] = \mathbb{E} \text{Exp}_{z \sim \mathcal{D}} [\mathbf{v}_t] = \nabla f_t(\mathbf{w}_t)$

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

# 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  $\nabla f(\mathbf{w}_t) = \mathbb{E} \text{Exp}_{z \sim \mathcal{D}} [\nabla \ell(\mathbf{w}_t, z)] = \mathbb{E} \text{Exp}_{z \sim \mathcal{D}} [\mathbf{v}_t] = \nabla f_t(\mathbf{w}_t)$

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