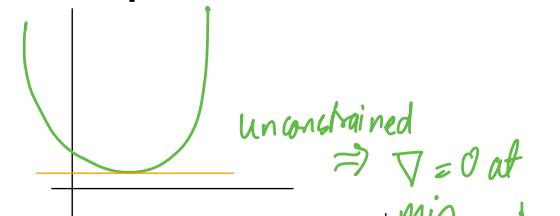


The Gradient Descent (GD) Algorithm

Outline:

$$C = \mathbb{R}^p$$

We will first consider the general form of the **basic** (non-stochastic) GD algorithm for **unconstrained** convex optimization of Lipschitz functions (not specifically in the learning context).



Then, we will consider a more advanced setting:

- **Constrained** optimization: $C \subset \mathbb{R}^p$



Required :

Solve $\min_{w \in C} f(w)$

where f is a convex, ρ -Lipschitz function (for some $\rho > 0$)

In this case, $w^* = \arg \min_{w \in C} f(w)$ is called the **constrained minimizer** of f over C .

The Gradient Descent (GD) Algorithm

Outline:

$$C = \mathbb{R}^p$$

We will first consider the general form of the **basic** (non-stochastic) GD algorithm for **unconstrained** convex optimization of Lipschitz functions (not specifically in the learning context).

Then, we will consider a more advanced setting:

- **Constrained** optimization: $C \subset \mathbb{R}^p$
- The gradient descent (GD) algorithm is an iterative algorithm for solving the above optimization problem (for both the unconstrained and constrained settings).
- As the number of iterations increases, the output of the GD algorithm converges to the minimizer w^* of f .
- Next, we will consider **Stochastic** GD (SGD) algorithm for minimizing the loss in bounded-convex-Lipschitz **learning** problems.

Basic GD algorithm (non-stochastic, unconstrained optimization)

- **Inputs:**

- A **convex**, ρ -**Lipschitz** function: $f : \mathbb{R}^p \rightarrow \mathbb{R}$

we will assume that f is **differentiable**, i.e., the gradient $\nabla f(\mathbf{w})$ exists for all $\mathbf{w} \in \mathbb{R}^p$.

- Number of iterations: T

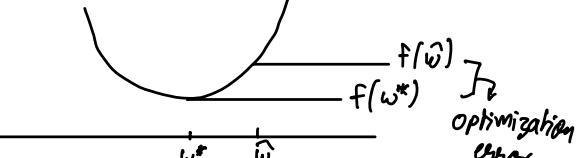
- Set of scalars: $\{\alpha_t : t = 1, \dots, T - 1\}$ (α_t called the **step size** or **learning rate**).

- **Output:**

A vector $\hat{\mathbf{w}}$: an approximate version of \mathbf{w}^* where $\mathbf{w}^* \in \arg \min_{\mathbf{w} \in \mathbb{R}^p} f(\mathbf{w})$

The quantity $f(\hat{\mathbf{w}}) - f(\mathbf{w}^*)$ is called the **optimization error**.

It measures the quality of the approximate solution $\hat{\mathbf{w}}$.



NOT $\hat{w} - w^*$

Basic GD algorithm (non-stochastic, unconstrained optimization)

- **Inputs:** $f : \mathbb{R}^p \rightarrow \mathbb{R}$, T , $\{\alpha_t : t = 1, \dots, T - 1\}$

1. Initialization: Choose an initial point $\mathbf{w}_1 = \mathbf{0}$ (all-zero vector in \mathbb{R}^p)

2. FOR $t = 1, \dots, T - 1$

Take a GD step: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha_t \nabla f(\mathbf{w}_t)$

3. Return $\hat{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$

Averaging helps smooth out gradient because of general-ness of our assumption gradient gives us dir of largest growth locally.

Theorem: (Convergence of Basic GD Algorithm)

In general,

$$\|\mathbf{w}^* - \mathbf{w}_t\| \leq M$$

w_1 is initialized w. our case = 0.

Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be a convex, ρ -Lipschitz function. Let $\mathbf{w}^* \in \arg \min_{\mathbf{w} \in \mathbb{R}^p} f(\mathbf{w})$

and suppose that $\|\mathbf{w}^*\| \leq M$. If we run the GD algorithm above for T steps

with $\alpha_t = \alpha = \frac{M}{\rho\sqrt{T}}$, $\forall t$. Then,

$$f(\hat{\mathbf{w}}) - f(\mathbf{w}^*) \leq \frac{M\rho}{\sqrt{T}}$$

optimization error

This is optimal convergence rate for this type of objective functions

error bound scales as $\frac{1}{\sqrt{T}}$
(goes to 0 as $T \uparrow \infty$)

OPTIMAL WRT First-order methods (using gradients).

Basic GD algorithm: Proof of the main theorem

Theorem: (Convergence of Basic GD Algorithm)

Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be a convex, ρ -Lipschitz function. Let $\mathbf{w}^* \in \arg \min_{\mathbf{w} \in \mathbb{R}^p} f(\mathbf{w})$ and suppose that $\|\mathbf{w}^*\| \leq M$. If we run the GD algorithm above for T steps with $\alpha_t = \alpha = M / (\rho \sqrt{T})$. Then, $f(\hat{\mathbf{w}}) - f(\mathbf{w}^*) \leq M \rho / \sqrt{T}$

First, let's define $\psi_t = \|\mathbf{w}_t - \mathbf{w}^*\|^2$, $\forall t$ (ψ_t usually called the potential).

We prove the theorem using the following claims:

$$\text{Claim 1: } \forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$$

The claims involve upper and lower bounds on the same quantity $\langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle$

$$\text{Claim 2: } \forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \geq f(\mathbf{w}_t) - f(\mathbf{w}^*)$$

Claim 1 gives an upper bound on $\langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle$ in terms of the difference of the potential functions at consecutive iterations $\psi_t - \psi_{t+1}$, step size α , and the Lipschitz constant ρ .

Basic GD algorithm: Proof of the main theorem

Theorem: (Convergence of Basic GD Algorithm)

Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be a convex, ρ -Lipschitz function. Let $\mathbf{w}^* \in \arg \min_{\mathbf{w} \in \mathbb{R}^p} f(\mathbf{w})$ and suppose that $\|\mathbf{w}^*\| \leq M$. If we run the GD algorithm above for T steps with $\alpha_t = \alpha = M / (\rho \sqrt{T})$. Then, $f(\hat{\mathbf{w}}) - f(\mathbf{w}^*) \leq M \rho / \sqrt{T}$

First, let's define $\psi_t = \|\mathbf{w}_t - \mathbf{w}^*\|^2$, $\forall t$ (ψ_t usually called the potential).

We prove the theorem using the following claims:

$$\text{Claim 1: } \forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$$

$$\text{Claim 2: } \forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \geq f(\mathbf{w}_t) - f(\mathbf{w}^*)$$

Claim 2 gives a lower bound on $\langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle$ in terms of the optimization error at iteration t : $f(\mathbf{w}_t) - f(\mathbf{w}^*)$.

Basic GD algorithm: Proof of the main theorem

Theorem: (Convergence of Basic GD Algorithm)

Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be a convex, ρ -Lipschitz function. Let $\mathbf{w}^* \in \arg \min_{\mathbf{w} \in \mathbb{R}^p} f(\mathbf{w})$ and suppose that $\|\mathbf{w}^*\| \leq M$. If we run the GD algorithm above for T steps with $\alpha_t = \alpha = M / (\rho \sqrt{T})$. Then, $f(\hat{\mathbf{w}}) - f(\mathbf{w}^*) \leq M \rho / \sqrt{T}$

First, let's define $\psi_t = \|\mathbf{w}_t - \mathbf{w}^*\|^2$, $\forall t$ (ψ_t usually called the potential).

We prove the theorem using the following claims:

$$\text{Claim 1: } \forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$$

$$\text{Claim 2: } \forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \geq f(\mathbf{w}_t) - f(\mathbf{w}^*)$$

The quantity $\langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle$ can be viewed as a "proxy" quantity: combining the two claims we can upper bound the optimization error at iteration t in terms of the difference of potential functions, step size, and Lipschitz constant: $f(\mathbf{w}_t) - f(\mathbf{w}^*) \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$

Basic GD algorithm: Proof of the main theorem

Theorem: (Convergence of Basic GD Algorithm)

Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be a convex, ρ -Lipschitz function. Let $\mathbf{w}^* \in \arg \min_{\mathbf{w} \in \mathbb{R}^p} f(\mathbf{w})$ and suppose that $\|\mathbf{w}^*\| \leq M$. If we run the GD algorithm above for T steps with $\alpha_t = \alpha = M / (\rho \sqrt{T})$. Then, $f(\hat{\mathbf{w}}) - f(\mathbf{w}^*) \leq M \rho / \sqrt{T}$

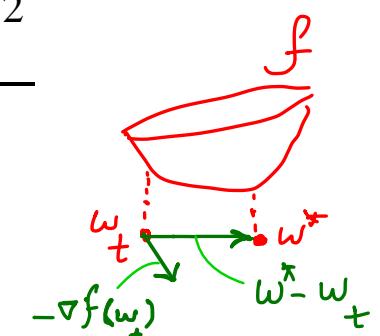
First, let's define $\psi_t = \|\mathbf{w}_t - \mathbf{w}^*\|^2$, $\forall t$ (ψ_t usually called the potential).

We prove the theorem using the following claims:

$$\text{Claim 1: } \forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$$

$$\text{Claim 2: } \forall t, \underbrace{\langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle}_{= \langle -\nabla f(\mathbf{w}_t), \mathbf{w}^* - \mathbf{w}_t \rangle} \geq f(\mathbf{w}_t) - f(\mathbf{w}^*)$$

$$= \langle -\nabla f(\mathbf{w}_t), \mathbf{w}^* - \mathbf{w}_t \rangle$$



For now, suppose the two claims are true.

(We will first show how to use the two claims to prove the theorem, then we will prove the two claims.)

Basic GD algorithm: Proof of the main theorem

Claim 1: $\forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \leq \frac{\Psi_t - \Psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$

Claim 2: $\forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \geq f(\mathbf{w}_t) - f(\mathbf{w}^*)$

Combining the two claims: $\forall t, f(\mathbf{w}_t) - f(\mathbf{w}^*) \leq \frac{\Psi_t - \Psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$ (#)

Note that our goal is to show Goal is to bound $\hat{\mathbf{w}}$ & not \mathbf{w}_t

$$f(\hat{\mathbf{w}}) - f(\mathbf{w}^*) \leq \frac{M\rho}{\sqrt{T}}$$

$\hat{\mathbf{w}}$ is the output of the algorithm = $\frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$

Let's see how can we show that given #

Basic GD algorithm: Proof of the main theorem

Claim 1: $\forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$

Claim 2: $\forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \geq f(\mathbf{w}_t) - f(\mathbf{w}^*)$

Combining the two claims: $\forall t, f(\mathbf{w}_t) - f(\mathbf{w}^*) \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$ (#)

Given those claims are true, now observe

$$f(\hat{\mathbf{w}}) - f(\mathbf{w}^*) \leq \frac{1}{T} \sum_{t=1}^T (f(\mathbf{w}_t) - f(\mathbf{w}^*))$$

By convexity of f and the definition of $\hat{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$

$$f\left(\frac{1}{T} \sum_{t=1}^T \mathbf{w}_t\right) \leq \frac{1}{T} \sum_{t=1}^T f(\mathbf{w}_t)$$

Basic GD algorithm: Proof of the main theorem

Claim 1: $\forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$

Claim 2: $\forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \geq f(\mathbf{w}_t) - f(\mathbf{w}^*)$

Combining the two claims: $\forall t, f(\mathbf{w}_t) - f(\mathbf{w}^*) \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$ (#)

Given those claims are true, now observe

$$f(\hat{\mathbf{w}}) - f(\mathbf{w}^*) \leq \frac{1}{T} \sum_{t=1}^T (f(\mathbf{w}_t) - f(\mathbf{w}^*))$$

By convexity of f and the definition of $\hat{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$

Recall: A fn $f: \mathbb{R}^P \rightarrow \mathbb{R}$ is convex if $\forall \lambda \in (0, 1)$
 $\forall u, v \in \mathbb{R}^P$ we have $f(\lambda u + (1-\lambda)v) \leq \lambda f(u) + (1-\lambda)f(v)$

In fact, this definition of convexity is equivalent to the following definition:

A fn $f: \mathbb{R}^P \rightarrow \mathbb{R}$ is convex if for any integer $K \geq 2$, any non-negative $\lambda_1, \dots, \lambda_K$ s.t. $\sum_{j=1}^K \lambda_j = 1$, and any $v_1, \dots, v_K \in \mathbb{R}^P$, we have $f\left(\sum_{j=1}^K \lambda_j v_j\right) \leq \sum_{j=1}^K \lambda_j f(v_j)$

$$w \in \mathbb{R}^d : \|w\| \leq 1$$

$$M=2$$

$$\rho = \underline{2(20)(10)}$$

$$(\langle w, x \rangle - y)^2$$

$$\nabla_w (\quad) = 2(\langle w, x \rangle - y)x$$

$$|\langle w, x \rangle - y| \leq 20$$

$$\|\cdot\| = 2 |(\langle w, x \rangle - y)| \|x\| \leq 10$$

$$\langle w, x \rangle \leq 10$$

Basic GD algorithm: Proof of the main theorem

Claim 1: $\forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$

Claim 2: $\forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \geq f(\mathbf{w}_t) - f(\mathbf{w}^*)$

Combining the two claims: $\forall t, f(\mathbf{w}_t) - f(\mathbf{w}^*) \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$ (#)

Given those claims are true, now observe

$$f(\hat{\mathbf{w}}) - f(\mathbf{w}^*) \leq \frac{1}{T} \sum_{t=1}^T (f(\mathbf{w}_t) - f(\mathbf{w}^*))$$

By convexity of f and the definition of $\hat{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$

Let $K = T$, and let $\lambda_1 = \lambda_2 = \dots = \lambda_T = \frac{1}{T}$ (note $\sum_{j=1}^T \lambda_j = 1$).

consider the parameter vectors $w_1, \dots, w_T \in \mathbb{R}^p$.

since f is convex, then

$$f\left(\sum_{j=1}^T \lambda_j w_j\right) \leq \sum_{j=1}^T \lambda_j f(w_j) \quad \text{i.e.,} \quad f\left(\frac{1}{T} \sum_{j=1}^T w_j\right) \leq \sum_{j=1}^T \frac{1}{T} f(w_j) = \frac{1}{T} \sum_{j=1}^T f(w_j)$$

Basic GD algorithm: Proof of the main theorem

Claim 1: $\forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$

Claim 2: $\forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \geq f(\mathbf{w}_t) - f(\mathbf{w}^*)$

Combining the two claims: $\forall t, f(\mathbf{w}_t) - f(\mathbf{w}^*) \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$ (#)

Given those claims are true, now observe

$$\begin{aligned} f(\hat{\mathbf{w}}) - f(\mathbf{w}^*) &\leq \frac{1}{T} \sum_{t=1}^T (f(\mathbf{w}_t) - f(\mathbf{w}^*)) \\ &\leq \frac{1}{2\alpha T} \sum_{t=1}^T (\psi_t - \psi_{t+1}) + \frac{\alpha\rho^2}{2} \\ &= \frac{1}{2\alpha T} (\psi_1 - \psi_{T+1}) + \frac{\alpha\rho^2}{2} \\ &\leq \frac{\psi_1}{2\alpha T} + \frac{\alpha\rho^2}{2} \leq \frac{M^2}{2\alpha T} + \frac{\alpha\rho^2}{2} \end{aligned}$$

By convexity of f and the definition of $\hat{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$

Follows directly from (#)

Telescopic sum: terms cancel.

Since $\psi_1 = \|\mathbf{w}^*\|^2 \leq M^2$

note $\psi_1 = \|\mathbf{w}_1 - \mathbf{w}^*\|^2 = \|\mathbf{w}^*\|^2$
since $\mathbf{w}_1 = \mathbf{0}$

Basic GD algorithm: Proof of the main theorem

Claim 1: $\forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$

Claim 2: $\forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \geq f(\mathbf{w}_t) - f(\mathbf{w}^*)$

Combining the two claims: $\forall t, f(\mathbf{w}_t) - f(\mathbf{w}^*) \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$ (#)

Given those claims are true, now observe

$$\begin{aligned} f(\hat{\mathbf{w}}) - f(\mathbf{w}^*) &\leq \frac{1}{T} \sum_{t=1}^T (f(\mathbf{w}_t) - f(\mathbf{w}^*)) \\ &\stackrel{\text{Telescopic Sum, will cancel}}{\leq} \frac{1}{2\alpha T} \sum_{t=1}^T (\psi_t - \psi_{t+1}) + \frac{\alpha\rho^2}{2} \\ &= \frac{1}{2\alpha T} (\psi_1 - \psi_{T+1}) + \frac{\alpha\rho^2}{2} \\ &\stackrel{\text{bounded by } M^2.}{\leq} \frac{\psi_1}{2\alpha T} + \frac{\alpha\rho^2}{2} \leq \frac{M^2}{2\alpha T} + \frac{\alpha\rho^2}{2} = \frac{M\rho}{\sqrt{T}} \end{aligned}$$

minimizing α gives this bound

By convexity of f and the definition of $\hat{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$

Follows directly from (#)

Telescopic sum: terms cancel.

Since $\psi_1 = \|\mathbf{w}^*\|^2 \leq M^2$

Last equality: by substitution with $\alpha = M / (\rho\sqrt{T})$

Basic GD algorithm: Proof of the main theorem

$$\text{Claim 1: } \forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$$

Proof :

$$\begin{aligned}\psi_{t+1} &= \|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2 = \|\mathbf{w}_t - \alpha \nabla f(\mathbf{w}_t) - \mathbf{w}^*\|^2 \\ &= \left\langle \underbrace{\mathbf{w}_t - \mathbf{w}^*}_{A} - \underbrace{\alpha \nabla f(\mathbf{w}_t)}_{B}, \underbrace{\mathbf{w}_t - \mathbf{w}^*}_{A} - \underbrace{\alpha \nabla f(\mathbf{w}_t)}_{B} \right\rangle \quad (A-B)^2 \\ &= \|\mathbf{w}_t - \mathbf{w}^*\|^2 - 2\alpha \underbrace{\left\langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \right\rangle}_{\text{in the claim}} + \alpha^2 \|\nabla f(\mathbf{w}_t)\|^2 \\ &\leq \psi_t - 2\alpha \left\langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \right\rangle + \alpha^2 \rho^2\end{aligned}$$

since $\|\nabla f(\mathbf{w}_t)\| \leq \rho$ by
 ρ -Lipschitzness of f

By rearranging the terms, the proof is complete.

Basic GD algorithm: Proof of the main theorem

Claim 2: $\forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \geq f(\mathbf{w}_t) - f(\mathbf{w}^*)$

Proof :

By convexity of f (recall the second definition of convex functions):

$$\begin{aligned} f(\mathbf{w}^*) &\geq f(\mathbf{w}_t) + \langle \nabla f(\mathbf{w}_t), \mathbf{w}^* - \mathbf{w}_t \rangle \\ &= f(\mathbf{w}_t) - \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \end{aligned}$$

Note the change in the sign in the second term inside the inner-product. *(change sign)*

By rearranging the terms, the proof is complete.

We just proved:

Theorem: (Convergence of Basic GD Algorithm)

Let $f: \mathbb{R}^p \rightarrow \mathbb{R}$ be a convex, ρ -Lipschitz function. Let $\mathbf{w}^* \in \arg \min_{\mathbf{w} \in \mathbb{R}^p} f(\mathbf{w})$ and suppose that $\|\mathbf{w}^*\| \leq M$. If we run the GD algorithm above for T steps with $\alpha_t = \alpha = M / (\rho \sqrt{T})$. Then, $f(\hat{\mathbf{w}}) - f(\mathbf{w}^*) \leq M \rho / \sqrt{T}$

The Gradient Descent (GD) Algorithm

Outline:

We will first consider the general form of the **basic** (non-stochastic) GD algorithm for **unconstrained** convex optimization of Lipschitz functions (not specifically in the learning context).



Then, we will consider a more advanced setting:

- **Constrained** optimization: $C \subset \mathbb{R}^p$  **Now**
- Next, we will consider **Stochastic** GD (SGD) algorithm for minimizing the loss in bounded-convex-Lipschitz learning problems.

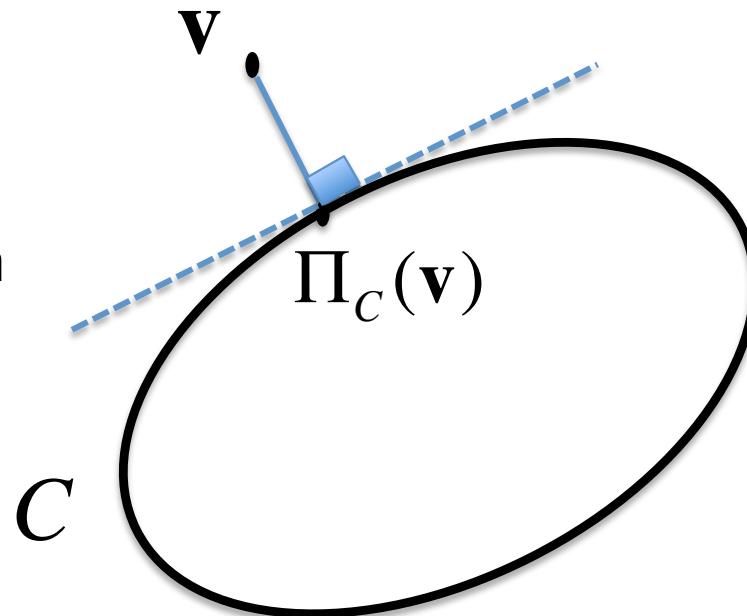
Euclidean projection

Definition: (Euclidean Projection)

Let $C \subseteq \mathbb{R}^p$ be a closed convex set. The Euclidean projection $\Pi_C : \mathbb{R}^p \rightarrow C$ is defined as:

$$\forall \mathbf{v} \in \mathbb{R}^p, \quad \Pi_C(\mathbf{v}) = \arg \min_{\mathbf{w} \in C} \|\mathbf{v} - \mathbf{w}\|$$

That is, $\Pi_C(\mathbf{v})$ is the “closest” point in C (w.r.t. the Euclidean distance) to \mathbf{v}



Euclidean projection

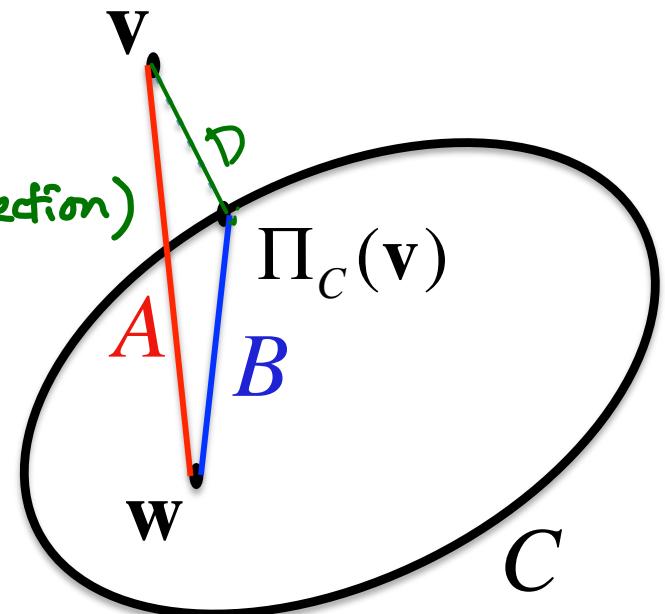
Fact:

$$\forall \mathbf{v} \in \mathbb{R}^p, \forall \mathbf{w} \in C, \text{ we have } \|\mathbf{v} - \mathbf{w}\| \geq \|\Pi_C(\mathbf{v}) - \mathbf{w}\|$$

Note that we are **not** aiming to prove that $A \geq D$

(which directly follows from the definition of Euclidean projection)

Before proving this fact, we will first see its useful application in the Gradient Descent algorithm.



$$A \geq B \text{ for any } \mathbf{w} \in C$$

*GD for **constrained** optimization: Projected GD*

- **Inputs:**

- A convex **ρ -Lipschitz** function: $f : \mathbb{R}^p \rightarrow \mathbb{R}$
- A convex constraint set: $C \subset \mathbb{R}^p$
- Number of iterations: T
- Set of scalars: $\{\alpha_t : t = 1, \dots, T - 1\}$ (**step size** or **learning rate**).

- **Output:**

An estimate $\hat{\mathbf{W}}$ of \mathbf{W}^* where $\mathbf{W}^* \in \arg \min_{\mathbf{w} \in C} f(\mathbf{w})$

GD for constrained optimization: Projected GD

- **Inputs:** $f, C, T, \{\alpha_t : t = 1, \dots, T - 1\}$
1. Initialization: $\mathbf{w}_1 = \mathbf{0}$ (assume all-zero vector is in C)
 2. **FOR** $t = 1, \dots, T - 1$
 - Take a **projected** GD step: $\mathbf{w}_{t+1} = \Pi_C(\mathbf{w}_t - \alpha_t \nabla f(\mathbf{w}_t))$
 3. **Return** $\hat{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$
- Note that if C is M -bounded, then this directly implies that $\|\mathbf{w}^* - \mathbf{w}_1\| \leq M$

Theorem: (Convergence of Projected GD Algorithm)

In general,

$$\|\mathbf{w}^* - \mathbf{w}_1\| \leq M$$

Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be a convex, ρ -Lipschitz function. Let $C \subset \mathbb{R}^p$ be a closed convex set. Let $\mathbf{w}^* \in \arg \min_{\mathbf{w} \in C} f(\mathbf{w})$ where $\|\mathbf{w}^*\| \leq M$. If we run the **projected** GD algorithm above for T steps with $\alpha_t = \alpha = \frac{M}{\rho\sqrt{T}}, \forall t$. Then, $f(\hat{\mathbf{w}}) - f(\mathbf{w}^*) \leq M\rho/\sqrt{T}$

GD for constrained optimization: Projected GD

Theorem: (Convergence of Projected GD Algorithm)

Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be a convex, ρ -Lipschitz function. Let $C \subset \mathbb{R}^p$ be a closed convex set. Let $\mathbf{w}^* \in \arg \min_{\mathbf{w} \in C} f(\mathbf{w})$ where $\|\mathbf{w}^*\| \leq M$. If we run the projected GD algorithm above for T steps with $\alpha_t = \alpha = \frac{M}{\rho\sqrt{T}}, \forall t$. Then, $f(\hat{\mathbf{w}}) - f(\mathbf{w}^*) \leq M\rho/\sqrt{T}$

As before, define $\psi_t = \|\mathbf{w}_t - \mathbf{w}^*\|^2, \forall t$

Recall the two claims used in the proof of the basic GD:

$$\textbf{Claim 1: } \forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$$

remains valid
(with an extra
step in the proof)

$$\textbf{Claim 2: } \forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \geq f(\mathbf{w}_t) - f(\mathbf{w}^*)$$

remains valid
(same proof:
convexity of f)

Projected GD algorithm: small modification in the proof

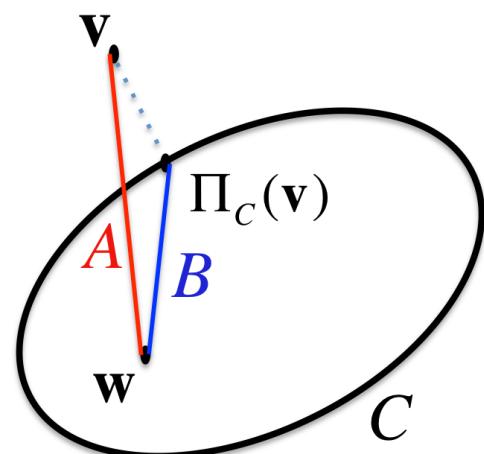
Claim 1: $\forall t, \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}^* \rangle \leq \frac{\psi_t - \psi_{t+1}}{2\alpha} + \frac{\alpha\rho^2}{2}$

Proof :

$$\begin{aligned} \psi_{t+1} &= \|\mathbf{w}_{t+1} - \mathbf{w}^*\|^2 = \|\Pi_C(\mathbf{w}_t - \alpha \nabla f(\mathbf{w}_t)) - \mathbf{w}^*\|^2 \\ &\stackrel{\substack{\text{upper bound since} \\ \text{projection on} \\ \text{convex maps is} \\ \text{a contractive mapping.}}}{\leq} \|\mathbf{w}_t - \alpha \nabla f(\mathbf{w}_t) - \mathbf{w}^*\|^2 \\ &\quad \vdots \\ &\quad \vdots \end{aligned}$$

Follows from the **Fact**:
the property of Euclidean
projections
(since $\mathbf{w}^* \in C$)

Continue the proof exactly as before.



$A \geq B$ for any $\mathbf{w} \in C$

Euclidean projection

Fact:

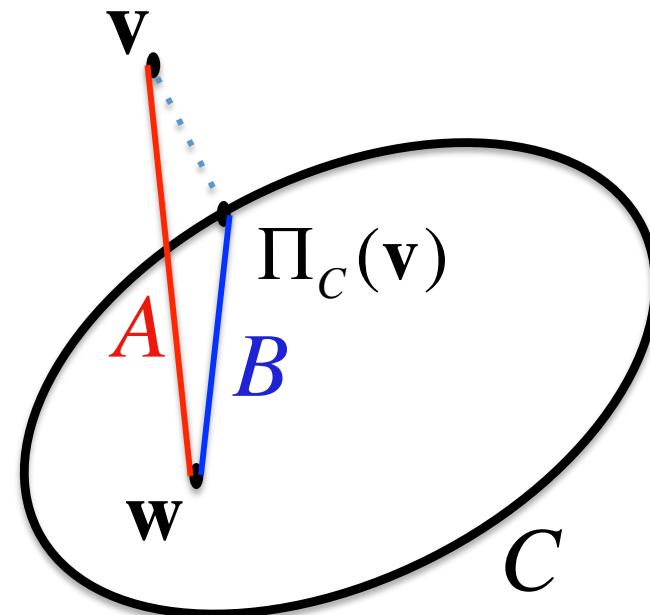
$$\forall \mathbf{v} \in \mathbb{R}^p, \forall \mathbf{w} \in C, \text{ we have } \|\mathbf{v} - \mathbf{w}\| \geq \|\Pi_C(\mathbf{v}) - \mathbf{w}\|$$

Proof:

$$\begin{aligned} \|\mathbf{v} - \mathbf{w}\|^2 &= \|\mathbf{v} - \Pi_C(\mathbf{v}) + \Pi_C(\mathbf{v}) - \mathbf{w}\|^2 \\ &= \|\Pi_C(\mathbf{v}) - \mathbf{w}\|^2 + 2\langle \mathbf{v} - \Pi_C(\mathbf{v}), \Pi_C(\mathbf{v}) - \mathbf{w} \rangle \\ &\quad + \|\mathbf{v} - \Pi_C(\mathbf{v})\|^2 \quad \leftarrow \text{Drop +ve term} \\ &\geq \|\Pi_C(\mathbf{v}) - \mathbf{w}\|^2 + 2\langle \mathbf{v} - \Pi_C(\mathbf{v}), \Pi_C(\mathbf{v}) - \mathbf{w} \rangle \end{aligned}$$

The following claim completes the proof

Claim: $\langle \mathbf{v} - \Pi_C(\mathbf{v}), \Pi_C(\mathbf{v}) - \mathbf{w} \rangle \geq 0$



$A \geq B$ for any $\mathbf{w} \in C$

Euclidean projection

Proof of the claim: $\langle \mathbf{v} - \Pi_C(\mathbf{v}), \Pi_C(\mathbf{v}) - \mathbf{w} \rangle \geq 0$

Since C is a convex set, then for any $\lambda \in (0, 1)$

$$(1-\lambda) \Pi_C(\mathbf{v}) + \lambda \mathbf{w} \in C$$

$\underbrace{\phantom{(1-\lambda) \Pi_C(\mathbf{v})}}_{\mathbf{u}_\lambda}$

By the definition of $\Pi_C(\mathbf{v})$

$$\begin{aligned} & \|\Pi_C(\mathbf{v}) - \mathbf{v}\|^2 \leq \|\mathbf{v} - \mathbf{u}_\lambda\|^2 \\ = & \|\Pi_C(\mathbf{v}) - \mathbf{v}\|^2 + 2\lambda \langle \mathbf{v} - \Pi_C(\mathbf{v}), \Pi_C(\mathbf{v}) - \mathbf{w} \rangle \\ & + \lambda^2 \|\Pi_C(\mathbf{v}) - \mathbf{w}\|^2 \end{aligned}$$

$$\text{Rearranging: } \langle \mathbf{v} - \Pi_C(\mathbf{v}), \Pi_C(\mathbf{v}) - \mathbf{w} \rangle \geq -\frac{\lambda}{2} \|\Pi_C(\mathbf{v}) - \mathbf{w}\|^2$$

Taking the limit as $\lambda \rightarrow 0$, we reach the desired result.

\angle is always $\geq 90^\circ$

