

COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

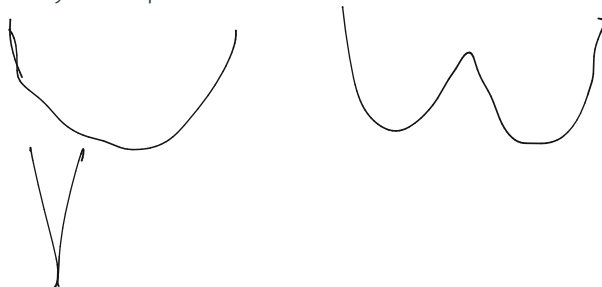
Cameron Musco

University of Massachusetts Amherst. Spring 2020.

Lecture 23

Last Class:

- Multivariable calculus review and gradient computation.
- Introduction to gradient descent. Motivation as a greedy algorithm.
- Conditions under which we will analyze gradient descent: convexity and Lipschitzness.



Last Class:

- Multivariable calculus review and gradient computation.
- Introduction to gradient descent. Motivation as a greedy algorithm.
- Conditions under which we will analyze gradient descent: convexity and Lipschitzness.

This Class:

- Analysis of gradient descent for Lipschitz, convex functions.
- Simple extension to projected gradient descent for constrained optimization.

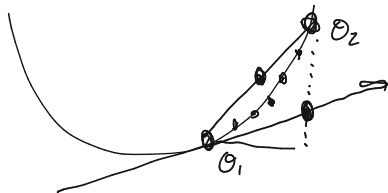
$$f(\theta) \quad \text{s.t.} \quad \theta \in S$$

Definition – Convex Function: A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if, for any $\vec{\theta}_1, \vec{\theta}_2 \in \mathbb{R}^d$ and $\lambda \in [0, 1]$:

$$(1 - \lambda) \cdot f(\vec{\theta}_1) + \lambda \cdot f(\vec{\theta}_2) \geq f\left((1 - \lambda) \cdot \vec{\theta}_1 + \lambda \cdot \vec{\theta}_2\right)$$

Corollary – Convex Function: A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if, for any $\vec{\theta}_1, \vec{\theta}_2 \in \mathbb{R}^d$ and $\lambda \in [0, 1]$:

$$\underline{f(\vec{\theta}_2) - f(\vec{\theta}_1)} \geq \underline{\vec{\nabla}f(\vec{\theta}_1)^T} (\vec{\theta}_2 - \vec{\theta}_1)$$



Definition – Convex Function: A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if, for any $\vec{\theta}_1, \vec{\theta}_2 \in \mathbb{R}^d$ and $\lambda \in [0, 1]$:

$$(1 - \lambda) \cdot f(\vec{\theta}_1) + \lambda \cdot f(\vec{\theta}_2) \geq f\left((1 - \lambda) \cdot \vec{\theta}_1 + \lambda \cdot \vec{\theta}_2\right)$$

Corollary – Convex Function: A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if, for any $\vec{\theta}_1, \vec{\theta}_2 \in \mathbb{R}^d$ and $\lambda \in [0, 1]$:

$$f(\vec{\theta}_2) - f(\vec{\theta}_1) \geq \vec{\nabla}f(\vec{\theta}_1)^T (\vec{\theta}_2 - \vec{\theta}_1)$$

Definition – Lipschitz Function: A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is G -Lipschitz if $\|\vec{\nabla}f(\vec{\theta})\|_2 \leq G$ for all $\vec{\theta}$.

Assume that:

$$\vec{\theta}_* = \arg \min_{\vec{\theta}} f(\vec{\theta})$$

- f is convex.
- f is G -Lipschitz.
- $\|\vec{\theta}_1 - \vec{\theta}_*\|_2 \leq R$ where $\vec{\theta}_1$ is the initialization point.

Gradient Descent

- Choose some initialization $\vec{\theta}_1$ and set $\eta = \frac{R}{G\sqrt{t}}$.
- For $i = 1, \dots, t-1$
 - $\vec{\theta}_{i+1} = \vec{\theta}_i - \eta \nabla f(\vec{\theta}_i)$
- Return $\hat{\theta} = \arg \min_{\vec{\theta}_1, \dots, \vec{\theta}_t} f(\vec{\theta}_i)$.

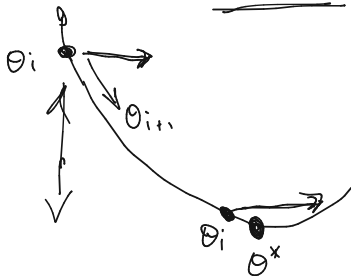
Theorem – GD on Convex Lipschitz Functions: For convex G -Lipschitz function f , GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon.$$

Theorem – GD on Convex Lipschitz Functions: For convex G -Lipschitz function f , GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon.$$

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$. Visually:



Theorem – GD on Convex Lipschitz Functions: For convex G -Lipschitz function f , GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon.$$

$$\|a-b\|_2^2 = \|a\|_2^2 + \|b\|_2^2 - 2a^T b$$

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$. Formally:

$$\begin{aligned} \|\theta_{i+1} - \theta_*\|_2^2 &= \|\theta_i - \eta \nabla f(\theta_i) - \theta_*\|_2^2 \\ &= \|\theta_i - \theta_*\|_2^2 + \underbrace{\|\eta \nabla f(\theta_i)\|_2^2}_{\leq \eta^2 G^2} - 2\eta \nabla f(\theta_i)^T (\theta_i - \theta_*) \end{aligned}$$

$$\|\theta_{i+1} - \theta_*\|_2^2 \leq \|\theta_i - \theta_*\|_2^2 + \eta^2 G^2 - 2\eta \nabla f(\theta_i)^T (\theta_i - \theta_*)$$

$$\nabla f(\theta_i)^T (\theta_i - \theta_*) \leq \frac{\|\theta_i - \theta_*\|_2^2 - \|\theta_{i+1} - \theta_*\|_2^2 + \eta^2 G^2}{2\eta}$$

Theorem – GD on Convex Lipschitz Functions: For convex G -Lipschitz function f , GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon.$$

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \underbrace{\frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta}} + \frac{\eta G^2}{2}.$

Step 1.1: $\underbrace{\vec{\nabla} f(\vec{\theta}_i)^T (\vec{\theta}_i - \vec{\theta}_*)}_{\leq 0} \leq \underbrace{\frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta}} + \frac{\eta G^2}{2}$

Theorem – GD on Convex Lipschitz Functions: For convex G -Lipschitz function f , GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

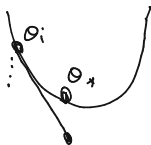
$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon.$$

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$.

Step 1.1: $\nabla f(\vec{\theta}_i)^T (\vec{\theta}_i - \vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2} \implies$ **Step 1.**

convexity

$$f(\theta_i) - f(\theta_*) \leq \nabla f(\theta_i)^T (\theta_i - \theta_*) \leq \text{RHS}$$



Theorem – GD on Convex Lipschitz Functions: For convex G -Lipschitz function f , GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon.$$

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$

Theorem – GD on Convex Lipschitz Functions: For convex G -Lipschitz function f , GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon.$$

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2} \implies$

telescoping sum

Step 2: $\frac{1}{t} \sum_{i=1}^t f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{R^2}{2\eta \cdot t} + \frac{\eta G^2}{2} \leq \epsilon$

want:
 $\min_{i:1 \dots t} f(\theta_i) - f(\theta_k) \leq \epsilon$

By Step 1:

$$\frac{1}{t} \sum f(\theta_i) - f(\theta_i) \leq \frac{mG^2}{2} + \frac{1}{t} \sum \left[\frac{\|\theta_i - \theta_*\|^2 - \|\theta_{i+1} - \theta_*\|^2}{2m} \right]$$

$$\underline{\|\theta_1 - \theta_*\|^2} - \underline{\|\theta_2 - \theta_*\|^2} + \underline{\|\theta_2 - \theta_*\|^2} - \underline{\|\theta_3 - \theta_*\|^2} + \dots - \underline{\|\theta_t - \theta_*\|^2} - \underline{\|\theta_{t+1} - \theta_*\|^2}$$

Theorem – GD on Convex Lipschitz Functions: For convex G -Lipschitz function f , GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of θ_* , outputs $\hat{\theta}$ satisfying:



$$f(\hat{\theta}) \leq f(\theta_*) + \epsilon.$$

distance you want covered



overshoot

Step 2: $\frac{1}{t} \sum_{i=1}^t f(\theta_i) - f(\theta_*) \leq \frac{R^2}{2\eta t} + \frac{\eta G^2}{2}$

$$\frac{nG^2}{2} + \frac{1}{t} \sum_{i=1}^t \frac{\|\theta_i - \theta_*\|_2^2 - \|\theta_{i+1} - \theta_*\|_2^2}{2n} = \frac{nG^2}{2} + \frac{\|\theta_1 - \theta_*\|_2^2 - \|\theta_{t+1} - \theta_*\|_2^2}{2nt}$$

$$\leq \frac{nG^2}{2} + \frac{R^2 - \|\theta_{t+1} - \theta_*\|_2^2}{2nt}$$

$$\leq \frac{nG^2}{2} + \frac{R^2}{2nt}$$

$n = \frac{R}{G\sqrt{t}}$

$$\Rightarrow \leq \frac{RG}{2\sqrt{t}} + \frac{R}{2\frac{R}{G\sqrt{t}}t} = \frac{RG}{2\sqrt{t}} + \frac{RG}{2\sqrt{t}} = \frac{RG}{\sqrt{t}} \leq \epsilon$$

CONSTRAINED CONVEX OPTIMIZATION

Often want to perform **convex optimization with convex constraints**.

$$\theta \in \mathbb{R}^d$$

$$\vec{\theta}^* = \arg \min_{\vec{\theta} \in \mathcal{S}} \underline{f(\vec{\theta})},$$

$$\|\theta\|_2^2 < 1$$

$\theta \in \text{subspace } V.$

where \mathcal{S} is a convex set.

CONSTRAINED CONVEX OPTIMIZATION

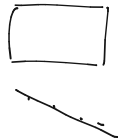
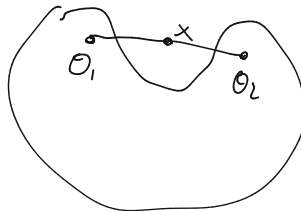
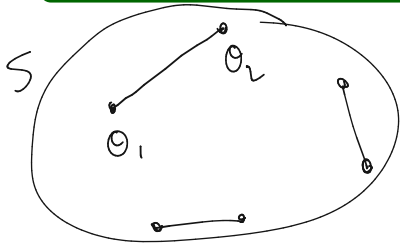
Often want to perform **convex optimization with convex constraints**.

$$\vec{\theta}^* = \arg \min_{\vec{\theta} \in \mathcal{S}} f(\vec{\theta}),$$

where \mathcal{S} is a **convex set**.

Definition – Convex Set: A set $\mathcal{S} \subseteq \mathbb{R}^d$ is convex if and only if, for any $\vec{\theta}_1, \vec{\theta}_2 \in \mathcal{S}$ and $\lambda \in [0, 1]$:

$$(1 - \lambda)\vec{\theta}_1 + \lambda \cdot \vec{\theta}_2 \in \mathcal{S}$$



CONSTRAINED CONVEX OPTIMIZATION

Often want to perform **convex optimization with convex constraints**.

$$\vec{\theta}^* = \arg \min_{\vec{\theta} \in \mathcal{S}} f(\vec{\theta}), \quad \textcircled{1} \quad \|a+b\| \leq \|a\| + \|b\|$$
$$\textcircled{2} \quad \|a-b\| \leq \|a\| + \|b\|$$
$$\|a-(-b)\| \leq \|a\| + \|(-b)\|$$

where \mathcal{S} is a **convex set**.

Definition – Convex Set: A set $\mathcal{S} \subseteq \mathbb{R}^d$ is convex if and only if, for any $\vec{\theta}_1, \vec{\theta}_2 \in \mathcal{S}$ and $\lambda \in [0, 1]$:

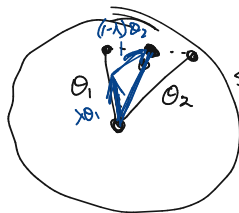
$$(1-\lambda)\vec{\theta}_1 + \lambda \cdot \vec{\theta}_2 \in \mathcal{S}$$

E.g. $\mathcal{S} = \{\vec{\theta} \in \mathbb{R}^d : \|\vec{\theta}\|_2 \leq 1\}$.

prove this is convex

$$\|\theta_1\| \leq 1, \quad \|\theta_2\| \leq 1$$

$$\|\lambda \theta_1 + (1-\lambda) \theta_2\| \leq 1$$

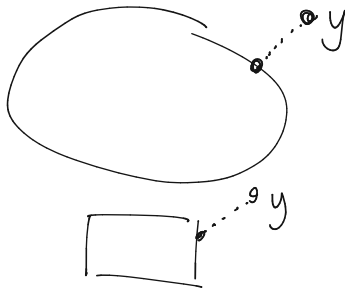


$$\|\lambda \theta_1 + (1-\lambda) \theta_2\|_2$$
$$\leq \lambda \|\theta_1\| + (1-\lambda) \|\theta_2\|_2$$
$$\leq \lambda \cdot 1 + (1-\lambda) \cdot 1$$
$$= 1$$

PROJECTED GRADIENT DESCENT

For any convex set let $P_S(\cdot)$ denote the projection function onto S .

$$\cdot \underline{P_S(\vec{y})} = \arg \min_{\vec{\theta} \in S} \|\vec{\theta} - \vec{y}\|_2.$$



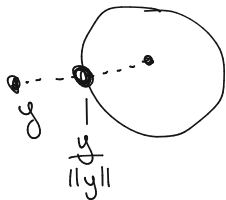
PROJECTED GRADIENT DESCENT

For any convex set let $P_S(\cdot)$ denote the projection function onto S .

- $P_S(\vec{y}) = \arg \min_{\vec{\theta} \in S} \|\vec{\theta} - \vec{y}\|_2$.

- For $S = \{\vec{\theta} \in \mathbb{R}^d : \|\vec{\theta}\|_2 \leq 1\}$ what is $P_S(\vec{y})$?

$$P_S(y) = \frac{y}{\|y\|_2}$$



PROJECTED GRADIENT DESCENT

For any convex set let $P_S(\cdot)$ denote the projection function onto S .

- $P_S(\vec{y}) = \arg \min_{\vec{\theta} \in S} \|\vec{\theta} - \vec{y}\|_2$.
- For $S = \{\vec{\theta} \in \mathbb{R}^d : \|\vec{\theta}\|_2 \leq 1\}$ what is $P_S(\vec{y})$?
- For S being a k dimensional subspace of \mathbb{R}^d , what is $P_S(\vec{y})$?



$$S = \{ \theta : \theta = Vc \}$$
$$P_S(y) = VV^T y$$

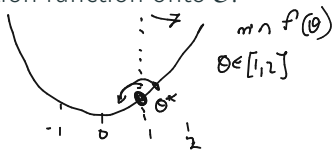
PROJECTED GRADIENT DESCENT

For any convex set let $P_S(\cdot)$ denote the projection function onto S .

- $P_S(\vec{y}) = \arg \min_{\vec{\theta} \in S} \|\vec{\theta} - \vec{y}\|_2$.

- For $S = \{\vec{\theta} \in \mathbb{R}^d : \|\vec{\theta}\|_2 \leq 1\}$ what is $P_S(\vec{y})$?

- For S being a k dimensional subspace of \mathbb{R}^d , what is $P_S(\vec{y})$?



Projected Gradient Descent

$$\|\theta\|_2^2 \leq 1 \text{ s.t. } \|\theta\|_2 \leq 1$$

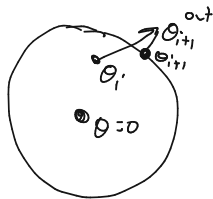
- Choose some initialization $\vec{\theta}_1$ and set $\eta = \frac{R}{G\sqrt{t}}$.

- For $i = 1, \dots, t-1$

- $\vec{\theta}_{i+1}^{(out)} = \vec{\theta}_i - \eta \cdot \vec{\nabla} f(\vec{\theta}_i)$

- $\vec{\theta}_{i+1} = P_S(\vec{\theta}_{i+1}^{(out)})$.

- Return $\hat{\theta} = \arg \min_{\vec{\theta}_i} f(\vec{\theta}_i)$.

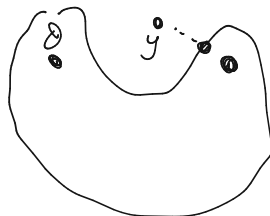
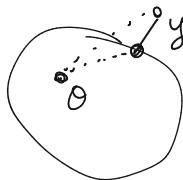


Projected gradient descent can be analyzed identically to gradient descent!

Projected gradient descent can be analyzed identically to gradient descent!

Theorem – Projection to a convex set: For any convex set $\mathcal{S} \subseteq \mathbb{R}^d$, $\vec{y} \in \mathbb{R}^d$, and $\vec{\theta} \in \mathcal{S}$,

$$\|P_{\mathcal{S}}(\vec{y}) - \vec{\theta}\|_2 \leq \|\vec{y} - \vec{\theta}\|_2.$$



Theorem – Projected GD: For convex G -Lipschitz function f , and convex set \mathcal{S} , Projected GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon = \min_{\vec{\theta} \in \mathcal{S}} f(\vec{\theta}) + \epsilon$$

Theorem – Projected GD: For convex G -Lipschitz function f , and convex set \mathcal{S} , Projected GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon = \min_{\vec{\theta} \in \mathcal{S}} f(\vec{\theta}) + \epsilon$$

Recall: $\vec{\theta}_{i+1}^{(out)} = \underline{\vec{\theta}_i} - \eta \cdot \vec{\nabla} f(\vec{\theta}_i)$ and $\vec{\theta}_{i+1} = \underline{P_{\mathcal{S}}(\vec{\theta}_{i+1}^{(out)})}$.

Theorem – Projected GD: For convex G -Lipschitz function f , and convex set \mathcal{S} , Projected GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon = \min_{\vec{\theta} \in \mathcal{S}} f(\vec{\theta}) + \epsilon$$

Recall: $\vec{\theta}_{i+1}^{(out)} = \vec{\theta}_i - \eta \cdot \vec{\nabla} f(\vec{\theta}_i)$ and $\vec{\theta}_{i+1} = P_{\mathcal{S}}(\vec{\theta}_{i+1}^{(out)})$.

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1}^{(out)} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$.

Theorem – Projected GD: For convex G -Lipschitz function f , and convex set S , Projected GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon = \min_{\vec{\theta} \in S} f(\vec{\theta}) + \epsilon$$

Recall: $\vec{\theta}_{i+1}^{(out)} = \vec{\theta}_i - \eta \cdot \vec{\nabla} f(\vec{\theta}_i)$ and $\vec{\theta}_{i+1} = P_S(\vec{\theta}_{i+1}^{(out)})$.

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1}^{(out)} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$.

[Step 1.a: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$.

$$\|\vec{\theta}_{i+1} - \vec{\theta}_*\| \leq \|\vec{\theta}_{i+1}^{out} - \vec{\theta}_*\|$$

[, ,]

Theorem – Projected GD: For convex G -Lipschitz function f , and convex set S , Projected GD run with $t \geq \frac{R^2 G^2}{\epsilon^2}$ iterations, $\eta = \frac{R}{G\sqrt{t}}$, and starting point within radius R of $\vec{\theta}_*$, outputs $\hat{\theta}$ satisfying:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon = \min_{\vec{\theta} \in S} f(\vec{\theta}) + \epsilon$$

Recall: $\vec{\theta}_{i+1}^{(out)} = \vec{\theta}_i - \eta \cdot \vec{\nabla} f(\vec{\theta}_i)$ and $\vec{\theta}_{i+1} = P_S(\vec{\theta}_{i+1}^{(out)})$.

Step 1: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1}^{(out)} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$.

Step 1.a: For all i , $f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{\|\vec{\theta}_i - \vec{\theta}_*\|_2^2 - \|\vec{\theta}_{i+1} - \vec{\theta}_*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$.

Step 2: $\frac{1}{t} \sum_{i=1}^t f(\vec{\theta}_i) - f(\vec{\theta}_*) \leq \frac{R^2}{2\eta \cdot t} + \frac{\eta G^2}{2} \implies$ Theorem.

Grad Descent:
 convex / Lipschitz
 constrained opt.
 over convex set
 via projected GD
 Project GD follow
 directly GD.