

MAT3007 Optimization

Lecture 19 Algorithms

Yuang Chen

School of Data Scienc
The Chinese University of Hong Kong, Shenzhen

July 17, 2025

Outline

- 1 Gradient Descent Review
- 2 Newton's Method
- 3 Projected Gradient Method
- 4 Why the Projected Gradient Method Works

Outline

- 1 Gradient Descent Review
- 2 Newton's Method
- 3 Projected Gradient Method
- 4 Why the Projected Gradient Method Works

Gradient Descent Algorithm

Start with any point \mathbf{x}^0 . Set $k = 0$ and stopping criterion $\epsilon > 0$

- ① Check $\|\nabla f(\mathbf{x}^k)\|$. If $\|\nabla f(\mathbf{x}^k)\| \leq \epsilon$, stop and output \mathbf{x}^k . Otherwise, continue to Step 2
- ② Let $\mathbf{d}^k = -\nabla f(\mathbf{x}^k)$
- ③ Use either exact line search or backtracking line search to find α_k
- ④ Let $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{d}^k$, let $k = k + 1$. Go back to step 1.

$$\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k - \alpha_k \nabla f(\mathbf{x}^k)$$

Step Size

- Constant step size: $\alpha_k = \bar{\alpha}$ for all k .
- Exact line search: choose α_k such that

$$\alpha_k = \arg \min_{\alpha} f(\mathbf{x}^k + \alpha \mathbf{d}^k)$$

- Backtracking / Armijo line search:
 - Let $\sigma, \gamma \in (0, 1)$ be given. Start with $\alpha = 1$.
 - If $f(\mathbf{x}^k + \alpha \mathbf{d}^k) \leq f(\mathbf{x}^k) + \gamma \alpha \cdot \nabla f(\mathbf{x}^k)^T \mathbf{d}^k$, choose $\alpha_k = \alpha$. Otherwise, set $\alpha = \sigma \alpha$ and repeat this step.

Gradient Descent Convergence Rates

Consider unconstrained problem $\min_{x \in \mathbb{R}^n} f(x)$.

- **Sublinear convergence:** If f is convex and ∇f is Lipschitz continuous, GD with constant step size has sublinear convergence (in function value) to the global minimizer.
- **Linear convergence:** If f is strongly convex and ∇f is Lipschitz continuous, GD with constant step size / exact / Armijo line search will have linear convergence to the unique global minimizer.
- **Global convergence:** No matter where GD starts, it will always converge to a local minimizer (when $f(\cdot)$ is convex, it converges to a global minimizer).

References and further reading

- S. Boyd and L. Vandenberghe (2004), *Convex Optimization*, Chapter 9
- T. Hastie, R. Tibshirani and J. Friedman (2009), *The Elements of Statistical Learning*, Chapters 10 and 16
- Y. Nesterov (1998), *Introductory Lectures on Convex Optimization: A Basic Course*, Chapter 2
- L. Vandenberghe, Lecture notes for EE 236C, UCLA

Outline

- 1 Gradient Descent Review
- 2 Newton's Method
- 3 Projected Gradient Method
- 4 Why the Projected Gradient Method Works

Intuition 1

- First order Taylor expansion for linearly approximating $f(x)$:

$$f(x + \Delta x) = f(x) + \nabla f(x)\Delta x + o(|\Delta x|).$$

$\overset{T}{\textcircled{1}}$ $\overset{\text{NN small}}{\textcircled{2}}$

- Take gradient on both sides:

$$\nabla f(x + \Delta x) = \nabla f(x) + \nabla^2 f(x)\Delta x + o(|\Delta x|). \quad \textcircled{3} = 0$$

$\textcircled{1}$ $\textcircled{2}$

- If we want $x^* = x + \Delta x$, then

$$\nabla f(x^*) = \nabla f(x + \Delta x) = \nabla f(x) + \nabla^2 f(x)\Delta x = 0$$

$$\Rightarrow x^* - x = \Delta x = -[\nabla^2 f(x)]^{-1}\nabla f(x) = -\frac{\nabla^2 f(x)}{\nabla^2 f(x)}$$

- Use this direction in each iteration

$$x^{k+1} = x^k - [\nabla^2 f(x^k)]^{-1}\nabla f(x^k)$$

Intuition 2

- Let x^k be the current iterate and consider a second-order Taylor approximation of $f(x)$ at x^k , i.e.

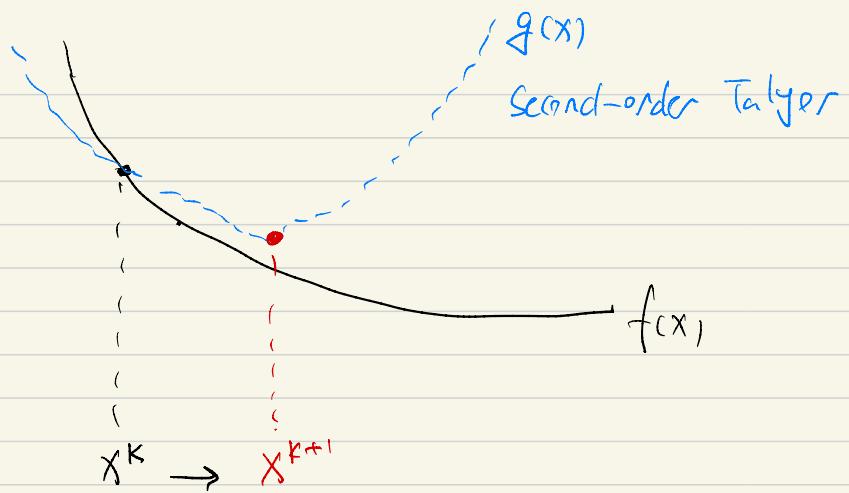
$$g(x) = f(x^k) + \nabla f(x^k)^\top (x - x^k) + \frac{1}{2}(x - x^k)^\top \nabla^2 f(x^k)(x - x^k).$$

- Choose the next iterate as the solution that minimizes the approximate function $g(x)$.
- Setting $\nabla g(x) = 0$, we get the linear system

$$\nabla f(x^k) + \nabla^2 f(x^k)(x - x^k) = 0.$$

- If the Hessian is non-singular, a solution to the above system is well-defined, and we set

$$x^{k+1} = x^k - [\nabla^2 f(x_k)]^{-1} \nabla f(x^k).$$



Newton's Method in High Dimensional Case

- Newton's step:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \underbrace{\frac{1}{(\nabla^2 f(\mathbf{x}^k))^{-1}}}_{\text{Step Size}} \nabla f(\mathbf{x}^k)$$

- Direction: $\mathbf{d}^k = -(\nabla^2 f(\mathbf{x}^k))^{-1} \nabla f(\mathbf{x}^k)$.
- Then Newton's step direction is a descent direction when f is convex.
- Newton's method may not converge unless the starting point is close.
- One way to help convergence is to use a step size parameter α_k in

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{d}^k$$

where $\mathbf{d}^k = -(\nabla^2 f(\mathbf{x}^k))^{-1} \nabla f(\mathbf{x}^k)$ is the Newton's step.

- One can use backtracking line search to determine α_k .
- In all above derivations, we assume that $\nabla^2 f(\mathbf{x}^k)$ is invertible.

$$d := -(\nabla^2 f(x))^{-1} \nabla f(x)$$

Prove d is a descent direction if f is convex.

Descent direction: $\{d : \nabla f(x)^T d < 0\}$

$$\begin{aligned}\nabla f(x)^T d &= -\underbrace{\nabla f(x)}_a^T \underbrace{(\nabla^2 f(x))^{-1}}_b \underbrace{\nabla f(x)}_a \\ &= -ba^2\end{aligned}$$

~~?~~ ~~0~~

If f is convex, then $\nabla^2 f(x) \succcurlyeq 0 \Rightarrow (\nabla^2 f(x))^{-1} \succcurlyeq 0$

Complete Procedure of Newton's Method

Start from a point \mathbf{x}^0 , set tolerance $\epsilon > 0$ and $k = 0$.

- ① Compute the Newton's step:

$$\mathbf{d}^k = -(\nabla^2 f(\mathbf{x}^k))^{-1} \nabla f(\mathbf{x}^k)$$

If $\|\nabla f(\mathbf{x}^k)\| < \epsilon$, then output \mathbf{x}^k , otherwise continue;

- ② Choose step size α_k by backtracking line search
- ③ $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{d}^k$, $k = k + 1$, go back to step 1

$$\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k - (\nabla^2 f(\mathbf{x}^k))^{-1} \nabla f(\mathbf{x}^k)$$

Linear Convergence vs Quadratic Convergence

Remember the gradient descent method can have linear convergence rate (strongly convex case):

$$|x^{k+1} - x^*| \leq \eta |x^k - x^*|.$$

Now, Newton's method has quadratic convergence rate:

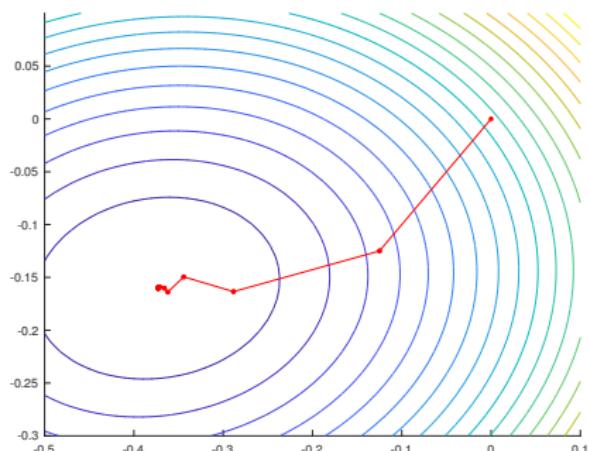
$$|x^{k+1} - x^*| \leq C |x^k - x^*|^2.$$

Convergence of Newton's Method

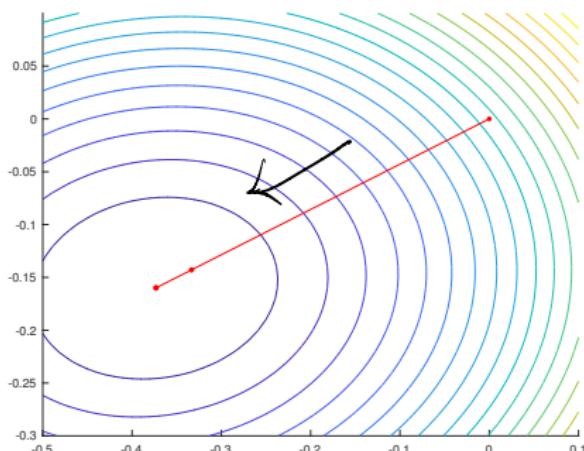
- Quadratic convergence in the neighborhood of a strict local minimum (under some conditions). Quadratic convergence is much faster than linear convergence in gradient descent.

Example: Minimize

$$f(x) = \exp(x_1 + x_2) + x_1^2 + 3x_2^2 - x_1x_2$$



(a) Gradient Descent



(b) Newton's Method

Newton's Method is Fast in Iteration Count, But ...

It requires computing the second-order derivative (Hessian) in each iteration.

- It might be computationally expensive to do so, especially if the second-order derivative doesn't have a closed-form (which is the case for many useful applications)
- It also requires a lot of space to store the Hessian matrix (If it is an n -dimensional problem, then we need n^2 space, comparing to n space required for gradient descent method).
- It requires much more matrix computation in each iteration
- If the Hessian is singular (or close to singular) at some iteration, we cannot proceed.
- Newton's method may diverge, even for convex functions.

if it starts from some point far away from x^* .

Quasi-Newton Methods

- Goal: Avoids computation of the Hessian and its inverse.
- Update iterate using

$\boxed{\nabla^2 f(x^k)}^{-1}$ *Newton*
↓

$$\rightarrow x^{k+1} = x^k - \alpha_k H_k \nabla f(x^k),$$

where α_k is the step size and H_k is an approximation to $[\nabla^2 f(x^k)]^{-1}$.

- Since H_{k+1}^{-1} approximates the Hessian,

$$H_{k+1}(\nabla f(x^{k+1}) - \nabla f(x^k)) = x^{k+1} - x^k.$$

- A widely used formula is the Broyden–Fletcher–Goldfarb–Shanno (BFGS) formula:

$$\curvearrowright H_{k+1} = H_k - \frac{d_k g_k^\top H_k + H_k g_k d_k^\top}{d_k^\top g_k} + \left(1 + \frac{g_k^\top H_k g_k}{d_k^\top g_k}\right) \frac{d_k d_k^\top}{d_k^\top g_k},$$

where $g_k = \nabla f(x^{k+1}) - \nabla f(x^k)$ and $d_k = x^{k+1} - x^k$.

Quiz

$$\text{max. } f(x) \rightarrow \text{-min. } -f(x)$$

What are the update equations in each iteration for gradient descent and Newton's method for a maximization problem?

$$\text{GD: } X^{k+1} \leftarrow X^k - \alpha \nabla f(X^k)$$

$$\text{Newton: } X^{k+1} \leftarrow X^k - [\nabla^2 f(X^k)]^{-1} \nabla f(X^k)$$

\downarrow
keeps
same

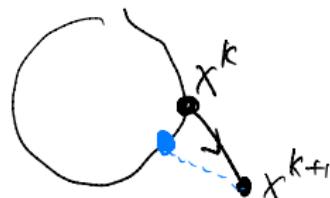
Outline

- 1 Gradient Descent Review
- 2 Newton's Method
- 3 Projected Gradient Method
- 4 Why the Projected Gradient Method Works

Constrained Optimization

We consider the following constrained optimization problem:

$$\begin{aligned} & \min_{\mathbf{x}} f(\mathbf{x}) \\ & \text{subject to } \mathbf{x} \in \Omega. \end{aligned}$$



where $\Omega \subset \mathbb{R}^n$ is a convex and closed set.

In the methods for unconstrained problems, the main idea is:

- At each \mathbf{x}^k , compute a descent direction \mathbf{d}^k .
- Then find an appropriate step size α_k and update $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{d}^k$.
- Both the gradient and Newton's method are based on this basic idea.
- \mathbf{x}^{k+1} can become infeasible. One solution is to project the point back to Ω .

Euclidean Projection / Orthogonal Projection

Definition: Euclidean / Orthogonal Projection

Let $\Omega \subset \mathbb{R}^n$ be a nonempty, closed, convex set. The (Euclidean / orthogonal) projection of x onto Ω is defined as the unique minimizer y^* of the constrained optimization problem:

$$\min_{\mathbf{y}} \quad \frac{1}{2} \| \mathbf{y} - \mathbf{x} \|^2 \quad \text{s.t.} \quad \mathbf{y} \in \Omega$$

and we write $y^* = \mathcal{P}_\Omega(x)$.

Observations:

- Interpretation: Projection is to find the closest point y^* in Ω to x .
- Hence, the projection $y^* = \mathcal{P}_\Omega(x)$ is the point in Ω that has the minimum distance to x .

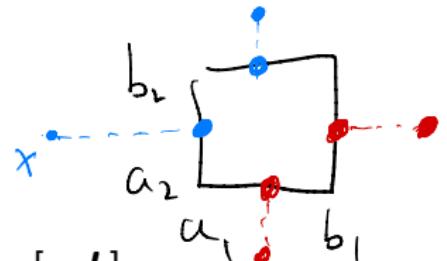
$$\boxed{\mathcal{P}_\Omega(x) = x \quad \text{if} \quad x \in \Omega}$$

Example I: Box Constraints

Suppose that Ω is given by box constraints:

$$\Omega = \{x \in \mathbb{R}^n : x_i \in [a_i, b_i], \forall i\} = [\mathbf{a}, \mathbf{b}],$$

where $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$, $\mathbf{a} \leq \mathbf{b}$, are given.



Euclidean Projection:

- The projection onto Ω can be computed as follows:

$$[\mathcal{P}_\Omega(x)]_i = \mathcal{P}_{[a_i, b_i]}(x_i) = \max \{\min\{x_i, b_i\}, a_i\}, \quad \forall i.$$



$$\mathcal{P}_{[a_i, b_i]}(x_i) = \begin{cases} a_i & \text{if } x_i \leq a_i \\ x_i & \text{if } x_i \in [a_i, b_i] \\ b_i & \text{if } x_i \geq b_i \end{cases}$$

Example II: Linear Constraints

We now consider the simple case where Ω consists of linear equality constraints:

$$\Omega = \{x : Ax = b\},$$

where $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$ are given.

Euclidean Projection:

- Suppose that A has full row rank ($m \leq n$), then it holds that:

$$\mathcal{P}_\Omega(x) = x - A^\top (A A^\top)^{-1} (Ax - b).$$

- This yields a special case when $\Omega = \{x : a^\top x = b\}$, i.e., projection onto a hyperplane. We have

$$\mathcal{P}_\Omega(x) = x - \frac{a^\top x - b}{\|a\|^2} a.$$

γ is given:

$$\begin{array}{ll} \min_{\gamma} & \frac{1}{2} \|\gamma - x\|^2 \\ \text{s.t.} & A\gamma = b \end{array} \quad \left. \begin{array}{l} \text{Convex} \\ \text{OPT} \end{array} \right\}$$

$$\text{Lagrangian: } L(\gamma, \mu) = \frac{1}{2} \|\gamma - x\|^2 + \mu^T (A\gamma - b)$$

KKT:

$$\text{main. } \nabla_{\gamma} L(\gamma, \mu) = \gamma - x + A^T \mu = 0$$

$$\text{Complementary slackness: } \mu^T (A\gamma - b) = 0$$

$$\text{Primal feasible: } A\gamma = b$$

$$\text{Dual feasible: } \text{only } \mu$$

$$A(\gamma - x + A^T \mu) = A(0)$$

$$A\gamma - Ax + AA^T \mu = 0$$

$$b - Ax + AA^T \mu = 0$$

$$\mu = (AA^T)^{-1}(Ax - b)$$

$$\gamma = x - A^T \mu = x - A^T (AA^T)^{-1}(Ax - b)$$

$$= P_{\Omega}(x)$$

Example III: Ball Constraints

Suppose that Ω is a Euclidean ball with radius $r > 0$ and center $\mathbf{m} \in \mathbb{R}^n$, i.e.:

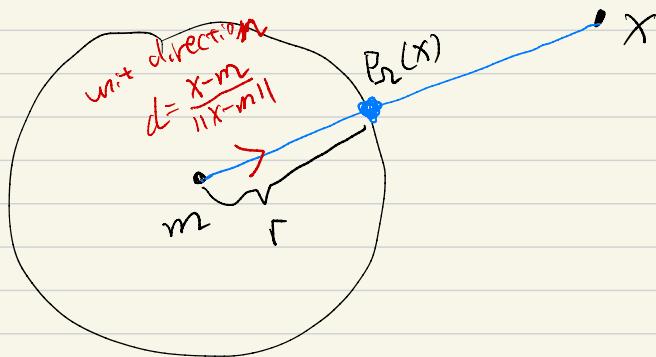
$$\Omega = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x} - \mathbf{m}\| \leq r\}.$$



Euclidean Projection:

- The projection onto Ω can be computed as follows:

$$\mathbf{y}^* = \mathcal{P}_\Omega(\mathbf{x}) = \begin{cases} \mathbf{x} & \text{if } \|\mathbf{x} - \mathbf{m}\| \leq r, \checkmark \\ \mathbf{m} + \frac{r}{\|\mathbf{x} - \mathbf{m}\|}(\mathbf{x} - \mathbf{m}) & \text{if } \|\mathbf{x} - \mathbf{m}\| > r. \end{cases}$$



$$P_r(x) = m + r \frac{x-m}{\|x-m\|} \quad \text{if } x \notin S_r$$

$\min_y \|y - x\|^2$
 $\text{s.t. } \|y - m\| \leq r$

Solve LKT

Projected Gradient Method: Basic Idea

In the gradient descent method, we perform a gradient step of the form:

$$\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k - \lambda_k \nabla f(\mathbf{x}^k),$$

where $\lambda_k > 0$ is a step size.

Motivation & Strategy:

- Setting $\mathbf{x}^{k+1} = \mathbf{x}^k - \lambda_k \nabla f(\mathbf{x}^k)$ might likely generate infeasible iterates: $\mathbf{x}^{k+1} \notin \Omega$.
- Idea: We project the step $\mathbf{x}^k - \lambda_k \nabla f(\mathbf{x}^k)$ back onto Ω :

$$\mathbf{x}^{k+1} = \underbrace{\mathcal{P}_\Omega(\mathbf{x}^k - \lambda_k \nabla f(\mathbf{x}^k))}_{\text{Projection}}.$$

Questions:

- How to choose $\lambda_k > 0$? If λ_k is determined by line search, every adjustment of λ_k requires re-evaluating the projection $\mathcal{P}_\Omega(\mathbf{x}^k - \lambda_k \nabla f(\mathbf{x}^k))$. This can be expensive.
- Can we guarantee descent?

Important Observation

Observation:

$$\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k + \alpha_k \mathbf{d}^k$$

$$\underline{\mathbf{x}^{k+1}} = \mathcal{P}_\Omega(\mathbf{x}^k - \lambda_k \nabla f(\mathbf{x}^k)) = \underline{\mathbf{x}^k} + \boxed{\mathcal{P}_\Omega(\mathbf{x}^k - \lambda_k \nabla f(\mathbf{x}^k)) - \underline{\mathbf{x}^k}}.$$

- This is close to our usual update form: $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{d}^k$.
- Setting $\mathbf{d}^k = \mathcal{P}_\Omega(\mathbf{x}^k - \lambda_k \nabla f(\mathbf{x}^k)) - \underline{\mathbf{x}^k}$, we can consider:

$$\rightarrow \mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{d}^k = (1 - \alpha_k) \underline{\mathbf{x}^k} + \alpha_k \mathcal{P}_\Omega(\mathbf{x}^k - \lambda_k \nabla f(\mathbf{x}^k)) \in \underline{\Omega}$$

If $\alpha_k \in [0, 1]$, then the convexity of Ω implies that \mathbf{x}^{k+1} will be *by* feasible if $\mathbf{x}^k \in \Omega$.

Why does / can this work? Is \mathbf{d}^k a descent direction?

Since convex set.

The Projected Gradient Method

Projected Gradient Method

- ① Initialization: Choose an initial point $\mathbf{x}^0 \in \Omega$.
- ② For $k = 0, 1, \dots$:
 - ① Select $\lambda_k > 0$ and compute $\nabla f(\mathbf{x}^k)$ and the search direction

$$\mathbf{d}^k = \mathcal{P}_\Omega(\mathbf{x}^k - \lambda_k \nabla f(\mathbf{x}^k)) - \mathbf{x}^k.$$

- ② Choose a step size $\alpha_k > 0$. $\alpha_k \in [0, 1]$
- ③ Set $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{d}^k$.
- ④ If $\|\mathbf{d}^{k+1}\| \leq \epsilon$, then STOP and \mathbf{x}^{k+1} is the output.

- In terms of convergence theorem, under appropriate conditions on the step sizes α_k and λ_k , projected gradient method has similar convergence guarantees to those of GD.

Outline

- 1 Gradient Descent Review
- 2 Newton's Method
- 3 Projected Gradient Method
- 4 Why the Projected Gradient Method Works

Optimization Problems with Convex Constraints

Recall the constrained optimization problem:

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{x} \in \Omega, \quad (1)$$

where $\Omega \subset \mathbb{R}^n$ is a convex and closed set.

- We can derive the following optimality condition:

Theorem: FONC for Problems with Convex Constraints

Let f be continuously differentiable on an open set that contains the convex and closed set $\Omega \subset \mathbb{R}^n$. Let $\mathbf{x}^* \in \Omega$ be a local minimizer of (1), then:

$$\nabla f(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*) \geq 0, \quad \forall \mathbf{x} \in \Omega. \quad (2)$$

- A point satisfying (2) is called a stationary point of problem (1).

Proof: Let x^* be a local optimal
 $\exists \varepsilon > 0$ s.t.

$$f(x^*) \leq f(x), \quad \forall x \in B_\varepsilon(x^*) \quad \text{(*)}$$

Let $x \in \mathcal{N}$ and $\lambda \in [0, 1]$

Define $x(\lambda) = \lambda x + (1-\lambda)x^*$ $\in \mathcal{N}$ because
 \mathcal{N} is convex set

Rewrite $x(\lambda) = x^* + \lambda(x - x^*)$

$$\begin{aligned} \|x(\lambda) - x^*\| &\leq \varepsilon \\ \Leftrightarrow \|\lambda(x - x^*)\| &\leq \varepsilon \\ \Leftrightarrow \lambda \|x - x^*\| &\leq \varepsilon \\ \Leftrightarrow \lambda &\leq \frac{\varepsilon}{\|x - x^*\|} \end{aligned}$$

when $\lambda \in [0, \frac{\varepsilon}{\|x - x^*\|}]$, $x(\lambda) \in B_\varepsilon(x^*)$

By (*), $f(x^*) \leq f(x(\lambda))$

$$\Leftrightarrow \underbrace{f(x(\lambda)) - f(x^*)}_{\lambda} \geq 0 \text{ if } \lambda \neq 0$$

Take limit of $\lambda \rightarrow 0$, we get

$$\Rightarrow f(x^*) (x - x^*) \geq 0, \quad \forall x \in \mathcal{N}$$

Projection Theorem

Projection Theorem

Let Ω be a nonempty, closed, and convex set. Then:

- 1 • A point \mathbf{y}^* is the projection of \mathbf{x} onto Ω , i.e., $\underline{\mathcal{P}_\Omega(\mathbf{x})}$, if and only if
$$(\mathbf{y}^* - \mathbf{x})^\top (\mathbf{y} - \mathbf{y}^*) \geq 0, \quad \forall \mathbf{y} \in \Omega.$$
- 2 • The mapping $\mathcal{P}_\Omega : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is Lipschitz continuous with constant $L = 1$, i.e.,

$$\|\mathcal{P}_\Omega(\mathbf{x}) - \mathcal{P}_\Omega(\mathbf{y})\| \leq \|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$$

- 3 • (FONC for Problem with Convex Constraints) For any $\lambda > 0$, the vector \mathbf{x}^* is a stationary point of (1) if and only if

$$\mathbf{x}^* - \mathcal{P}_\Omega(\mathbf{x}^* - \lambda \nabla f(\mathbf{x}^*)) = 0.$$

$$1. \underset{y}{\text{m.h.}} \frac{1}{2} \|y - x\|^2 = f(y)$$

s.t. $y \in \Omega$

Plug in FONC (2): $\Rightarrow f(y^*) (y - y^*) \geq 0, \forall y \in \Omega$

$$\Rightarrow \underline{(y^* - x) (y - y^*) \geq 0, \forall y \in \Omega}$$

2. From (1), $y^* = P_{\Omega}(x)$

$$(P_{\Omega}(x) - x)^T (z - P_{\Omega}(x)) \geq 0, \forall z \in \Omega$$

From (2), $y^* = P_{\Omega}(y)$

Pick
 $z = P_{\Omega}(y)$

$$(P_{\Omega}(y) - y)^T (z - P_{\Omega}(y)) \geq 0, \forall z \in \Omega$$

Pick
 $z = P_{\Omega}(x)$

$$\begin{cases} (P_{\Omega}(x) - x)^T (P_{\Omega}(y) - P_{\Omega}(x)) \geq 0 \\ (P_{\Omega}(y) - y)^T (P_{\Omega}(x) - P_{\Omega}(y)) \geq 0 \end{cases}$$

$$\Rightarrow \underbrace{(P_{\Omega}(x) - P_{\Omega}(y))^T}_{(P_{\Omega}(y) - P_{\Omega}(x))} \underbrace{(P_{\Omega}(y) - P_{\Omega}(x) + x - y)}_{(x-y)} \geq 0$$

$$\Rightarrow -\|P_{\Omega}(x) - P_{\Omega}(y)\|^2 + (P_{\Omega}(x) - P_{\Omega}(y))^T (x - y) \geq 0$$

$$\Rightarrow \|(P_{\mathcal{L}}(x) - P_{\mathcal{L}}(y))\|^2 \leq (P_{\mathcal{L}}(x) - P_{\mathcal{L}}(y))^T(x-y)$$

$$\Rightarrow \|(P_{\mathcal{L}}(x) - P_{\mathcal{L}}(y))\|^2 \leq \|(P_{\mathcal{L}}(x) - P_{\mathcal{L}}(y))\| \|(x-y)\|$$

$$\Rightarrow \|(P_{\mathcal{L}}(x) - P_{\mathcal{L}}(y))\| \leq \|x-y\|$$

3. $x^* = P_{\mathcal{L}}(x^* - \lambda \triangleright f(x^*))$

From 1. $\Leftrightarrow (x^* - (x^* - \lambda \triangleright f(x^*)))^T(x - x^*) \geq 0, \forall x \in \mathbb{R}$

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$
 Proj inside my Proj
 Proj in \mathcal{L}

$$\Leftrightarrow \lambda \triangleright f(x^*)^T(x - x^*) \geq 0 \quad \forall x \in \mathbb{R}$$

Since $\lambda > 0$

$$\Leftrightarrow \triangleright f(x^*)(x - x^*) \geq 0 \quad \forall x \in \mathbb{R}$$

Fo NC

Descent Direction

Lemma: Descent Direction

Let $\mathbf{x} \in \Omega$ and $\lambda > 0$ be given. If \mathbf{x} is not a stationary point of (1), then the direction

$$\mathbf{d} := \mathcal{P}_\Omega(\mathbf{x} - \lambda \nabla f(\mathbf{x})) - \mathbf{x}$$

is a descent direction and it holds that

$$\Rightarrow \nabla f(\mathbf{x})^\top \mathbf{d} \leq -\frac{1}{\lambda} \|\mathbf{d}\|^2 < 0.$$

\mathbf{d} is called the generalized gradient.

Proof:

$x \in S$

$$\nabla f(x)^T d$$

$$= \nabla f(x)^T (P_S(x - \lambda \nabla f(x)) - x)$$

$$= \nabla f(x)^T (P_S(x - \lambda \nabla f(x)) - P_S(x))$$

$$= -\frac{1}{\lambda} (x - \lambda \nabla f(x) - x)^T (P_S(x - \lambda \nabla f(x)) - P_S(x))$$

Cauchy-Schwarz $\leq -\frac{1}{\lambda} \|x - \lambda \nabla f(x) - x\| \cdot \|P_S(x - \lambda \nabla f(x)) - P_S(x)\|$

$$\leq -\frac{1}{\lambda} \|P_S(x - \lambda \nabla f(x)) - P_S(x)\| \cdot \|P_S(x - \lambda \nabla f(x)) - P_S(x)\|$$

$$= -\frac{1}{\lambda} \|P_S(x - \lambda \nabla f(x)) - P_S(x)\|^2$$

$$\leq -\frac{1}{\lambda} \|P_S(x - \lambda \nabla f(x)) - x\|^2$$

$$= -\frac{1}{\lambda} \|d\|^2 < 0$$