



Aalto University
School of Electrical
Engineering

Safety and Constrained Optimal Control

Gökhan Alcan

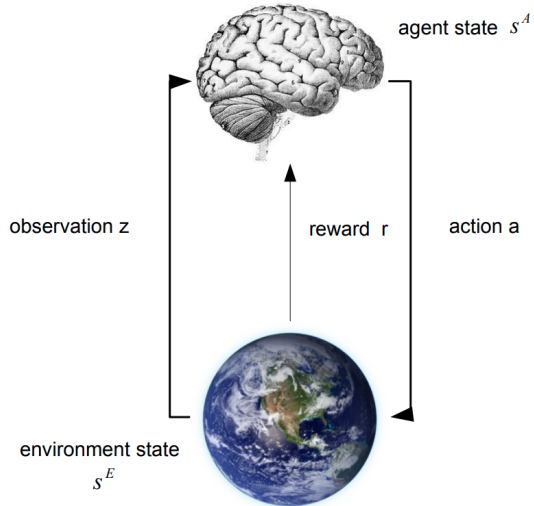
 Dept. of Electrical Engineering and Automation

 gokhan.alcan@aalto.fi

 www.gokhanalcan.com

November 9, 2021

Reinforcement Learning



Safety in Reinforcement Learning

- ▶ How would you define *safety* in RL?

Safety in Reinforcement Learning

- ▶ How would you define *safety* in RL?
- ▶ Safety in RL is an active research topic!

Safety in Reinforcement Learning

- ▶ How would you define *safety* in RL?
- ▶ Safety in RL is an active research topic!
- ▶ The agent is trained to *maximize the expected return* in a given task ...

Safety in Reinforcement Learning

- ▶ How would you define *safety* in RL?
- ▶ Safety in RL is an active research topic!
- ▶ The agent is trained to *maximize the expected return* in a given task *while not taking any action* that *gives damage* to the environment or itself during learning and/or deployment.

Safe Decision Making

From Control Theory Perspective

- ▶ Adaptive control

Safe Decision Making

From Control Theory Perspective

- ▶ Adaptive control
- ▶ Robust control

Safe Decision Making

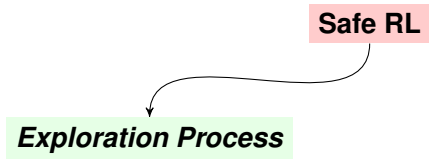
From Control Theory Perspective

- ▶ Adaptive control
- ▶ Robust control
- ▶ Robust model predictive control

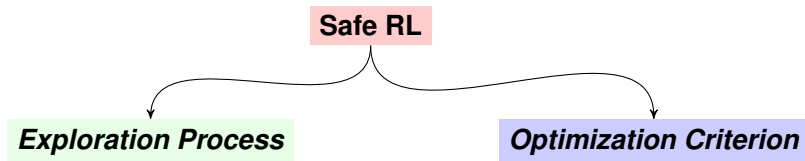
Safety in Reinforcement Learning

Safe RL

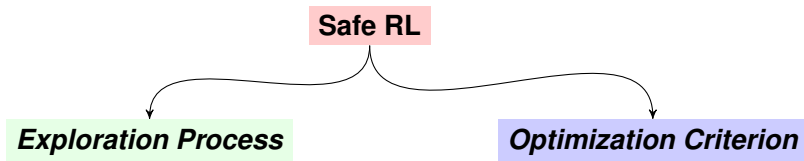
Safety in Reinforcement Learning



Safety in Reinforcement Learning

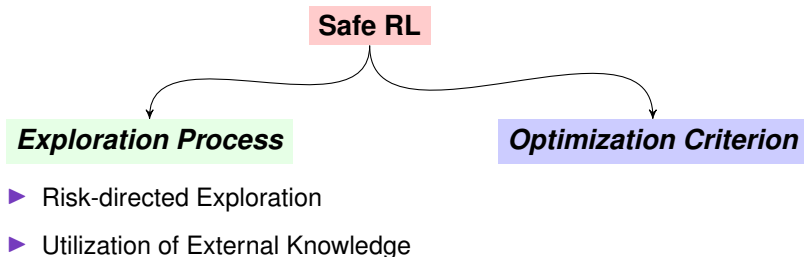


Safety in Reinforcement Learning

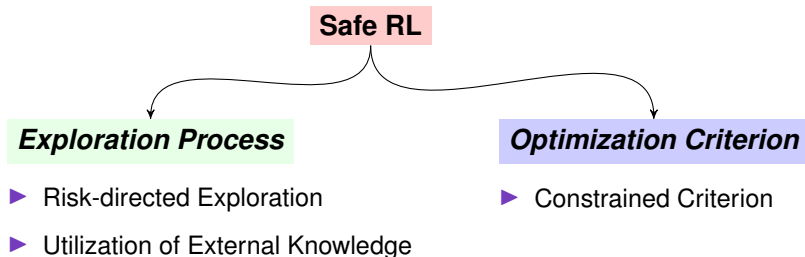


- ▶ Risk-directed Exploration

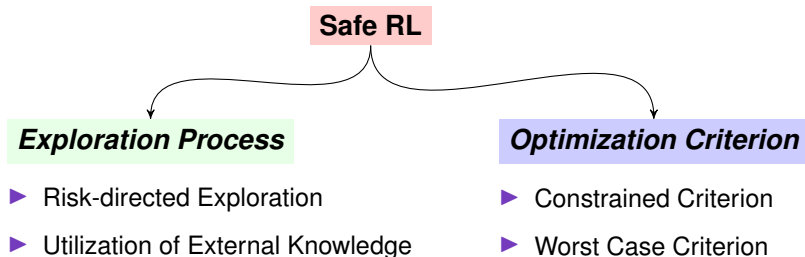
Safety in Reinforcement Learning



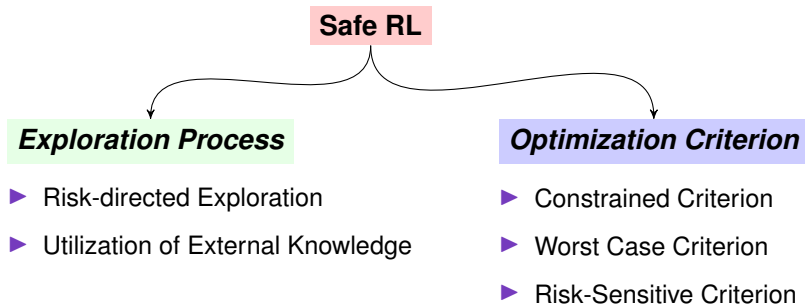
Safety in Reinforcement Learning



Safety in Reinforcement Learning

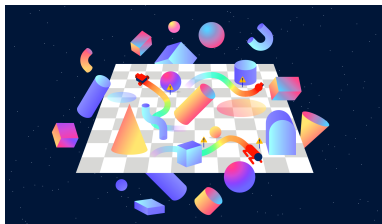


Safety in Reinforcement Learning



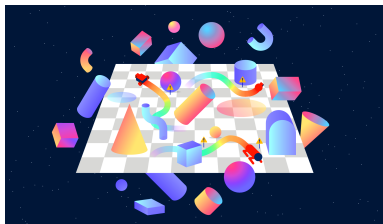
Safe Exploration

OpenAI Safety-Gym



Safe Exploration

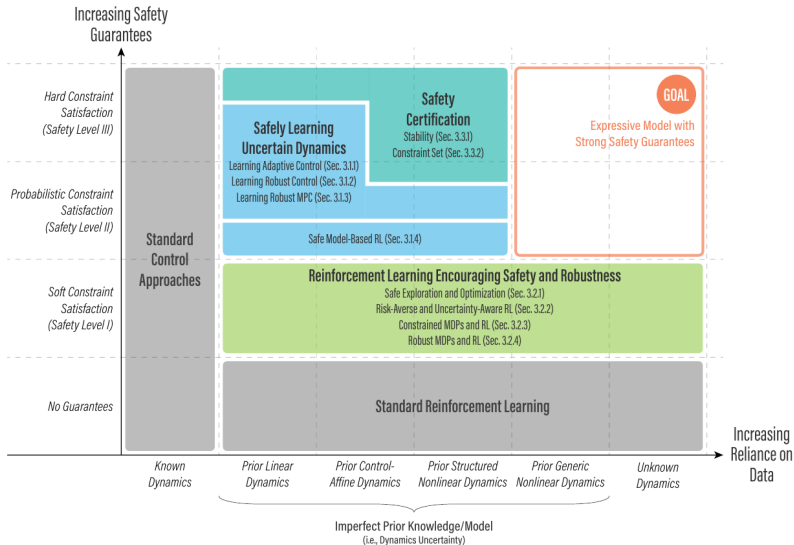
OpenAI Safety-Gym



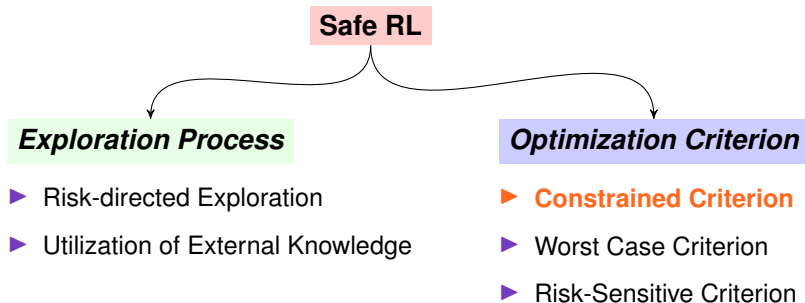
Some Methods

- ▶ Constrained Policy Optimization
- ▶ Proximal Policy Optimization
- ▶ Trust Region Policy Optimization
- ▶ PPO Lagrangian
- ▶ TRPO Lagrangian

Bridging Control Theory and RL



Safety in Reinforcement Learning



Constrained Optimal Control

Constrained Optimization

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{subject to} \quad \begin{cases} c_i(x) = 0, & i \in \mathcal{E} & \text{Equality Constraints} \\ c_i(x) \geq 0, & i \in \mathcal{I} & \text{Inequality Constraints} \end{cases}$$

Constrained Optimization

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{subject to} \quad \begin{cases} c_i(x) = 0, & i \in \mathcal{E} & \text{Equality Constraints} \\ c_i(x) \geq 0, & i \in \mathcal{I} & \text{Inequality Constraints} \end{cases}$$

Feasible Set:

$$\Omega = \{x \mid c_i(x) = 0, i \in \mathcal{E} \quad \text{and} \quad c_i(x) \geq 0, i \in \mathcal{I}\}$$

$$\implies \min_{x \in \Omega} f(x)$$

Constrained Optimization

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{subject to} \quad \begin{cases} c_i(x) = 0, & i \in \mathcal{E} & \text{Equality Constraints} \\ c_i(x) \geq 0, & i \in \mathcal{I} & \text{Inequality Constraints} \end{cases}$$

Feasible Set:

$$\Omega = \{x \mid c_i(x) = 0, i \in \mathcal{E} \quad \text{and} \quad c_i(x) \geq 0, i \in \mathcal{I}\}$$

$$\implies \min_{x \in \Omega} f(x)$$

Active Set:

$$\mathcal{A}(x) = \mathcal{E} \cup \{i \in \mathcal{I} \mid c_i(x) = 0\}$$

At a feasible point x , the inequality constraint $i \in \mathcal{I}$ is said to be **active** if $c_i(x) = 0$ and **inactive** if the strict inequality $c_i(x) > 0$ is satisfied.

Constrained Optimization

A Single Equality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad x_1^2 + x_2^2 - 2 = 0$$



Constrained Optimization

A Single Equality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad x_1^2 + x_2^2 - 2 = 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = x_1^2 + x_2^2 - 2$$

$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$



Constrained Optimization

A Single Equality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad x_1^2 + x_2^2 - 2 = 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = x_1^2 + x_2^2 - 2$$

$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

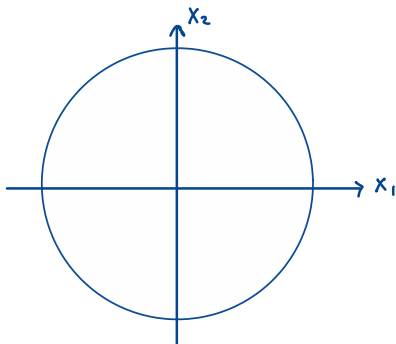
Q: What is feasible set?



Constrained Optimization

A Single Equality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad x_1^2 + x_2^2 - 2 = 0$$



$$f(x) = x_1 + x_2$$

$$c_1(x) = x_1^2 + x_2^2 - 2$$

$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

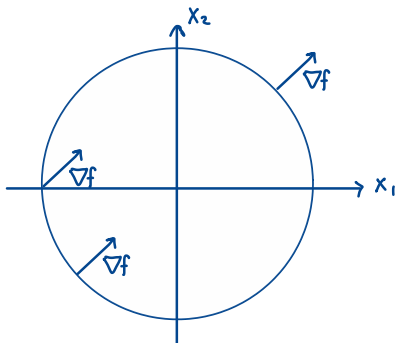
Q: What is feasible set?

A: Feasible set for this problem is a circle of radius $\sqrt{2}$ centered at origin.
(Just boundary, not interior)

Constrained Optimization

A Single Equality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad x_1^2 + x_2^2 - 2 = 0$$



$$f(x) = x_1 + x_2$$

$$c_1(x) = x_1^2 + x_2^2 - 2$$

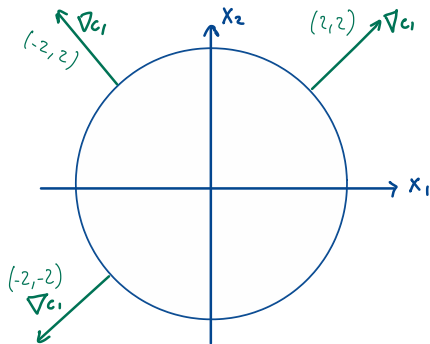
$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Constrained Optimization

A Single Equality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad x_1^2 + x_2^2 - 2 = 0$$



$$f(x) = x_1 + x_2$$

$$c_1(x) = x_1^2 + x_2^2 - 2$$

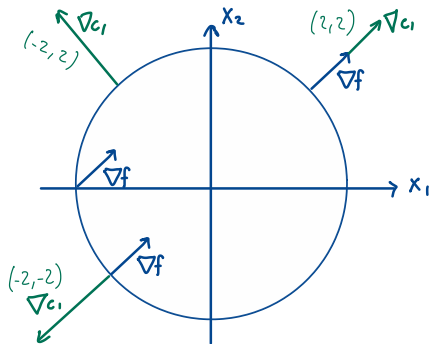
$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} 2x_1 \\ 2x_2 \end{bmatrix}$$

Constrained Optimization

A Single Equality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad x_1^2 + x_2^2 - 2 = 0$$



$$f(x) = x_1 + x_2$$

$$c_1(x) = x_1^2 + x_2^2 - 2$$

$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

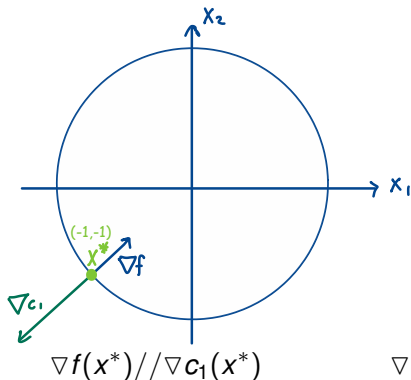
$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} 2x_1 \\ 2x_2 \end{bmatrix}$$

Q: What is the solution x^* ?

Constrained Optimization

A Single Equality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad x_1^2 + x_2^2 - 2 = 0$$



$$f(x) = x_1 + x_2$$

$$c_1(x) = x_1^2 + x_2^2 - 2$$

$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} 2x_1 \\ 2x_2 \end{bmatrix}$$

Q: What is the solution x^* ?

$$\mathbf{A: } x^* = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$

$$\nabla f(x^*) = \lambda_1^* \nabla c_1(x^*) \quad \lambda_1^* = -1/2$$

Constrained Optimization

A Single Equality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad x_1^2 + x_2^2 - 2 = 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = x_1^2 + x_2^2 - 2$$

$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

Let's introduce **Lagrangian function**

$$\mathcal{L}(x, \lambda_1) = f(x) - \lambda_1 c_1(x)$$

Constrained Optimization

A Single Equality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad x_1^2 + x_2^2 - 2 = 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = x_1^2 + x_2^2 - 2$$

$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

Let's introduce **Lagrangian function**

$$\mathcal{L}(x, \lambda_1) = f(x) - \lambda_1 c_1(x)$$

At solution x^* , there is a scalar λ_1^* such that $\nabla_x \mathcal{L}(x^*, \lambda_1^*) = 0$

Constrained Optimization

A Single Equality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad x_1^2 + x_2^2 - 2 = 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = x_1^2 + x_2^2 - 2$$

$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

Let's introduce **Lagrangian function**

$$\mathcal{L}(x, \lambda_1) = f(x) - \lambda_1 c_1(x)$$

At solution x^* , there is a scalar λ_1^* such that $\nabla_x \mathcal{L}(x^*, \lambda_1^*) = 0$

$$\nabla_x \mathcal{L}(x, \lambda_1) = \nabla f(x) - \lambda_1 \nabla c_1(x)$$

$$1 - 2\lambda_1^* x_1 = 0 \quad \text{and} \quad 1 - 2\lambda_1^* x_2 = 0$$

Constrained Optimization

A Single Equality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad x_1^2 + x_2^2 - 2 = 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = x_1^2 + x_2^2 - 2$$

$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

Let's introduce **Lagrangian function**

$$\mathcal{L}(x, \lambda_1) = f(x) - \lambda_1 c_1(x)$$

At solution x^* , there is a scalar λ_1^* such that $\nabla_x \mathcal{L}(x^*, \lambda_1^*) = 0$

$$\nabla_x \mathcal{L}(x, \lambda_1) = \nabla f(x) - \lambda_1 \nabla c_1(x)$$

$$1 - 2\lambda_1^* x_1 = 0 \quad \text{and} \quad 1 - 2\lambda_1^* x_2 = 0$$

Let's check our solution $x^* = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$, $\lambda_1^* = -1/2$

Constrained Optimization

A Single Equality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad x_1^2 + x_2^2 - 2 = 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = x_1^2 + x_2^2 - 2$$

$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

Let's introduce **Lagrangian function**

$$\mathcal{L}(x, \lambda_1) = f(x) - \lambda_1 c_1(x)$$

At solution x^* , there is a scalar λ_1^* such that $\nabla_x \mathcal{L}(x^*, \lambda_1^*) = 0$

$$\nabla_x \mathcal{L}(x, \lambda_1) = \nabla f(x) - \lambda_1 \nabla c_1(x)$$

$$1 - 2\lambda_1^* x_1 = 0 \quad \text{and} \quad 1 - 2\lambda_1^* x_2 = 0$$

Let's check our solution $x^* = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$, $\lambda_1^* = -1/2$

$$1 - 2(-1/2)(-1) = 0 \quad \text{and} \quad 1 - 2(-1/2)(-1) = 0 \quad \checkmark$$

Constrained Optimization

A Single Equality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad x_1^2 + x_2^2 - 2 = 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = x_1^2 + x_2^2 - 2$$

$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

Let's introduce **Lagrangian function**

$$\mathcal{L}(x, \lambda_1) = f(x) - \lambda_1 c_1(x)$$

At solution x^* , there is a scalar λ_1^* such that $\nabla_x \mathcal{L}(x^*, \lambda_1^*) = 0$

$$\nabla_x \mathcal{L}(x, \lambda_1) = \nabla f(x) - \lambda_1 \nabla c_1(x)$$

$$1 - 2\lambda_1^* x_1 = 0 \quad \text{and} \quad 1 - 2\lambda_1^* x_2 = 0$$

Q: What about $x = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$, $\lambda_1 = 1/2$?

Constrained Optimization

A Single Equality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad x_1^2 + x_2^2 - 2 = 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = x_1^2 + x_2^2 - 2$$

$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

Let's introduce **Lagrangian function**

$$\mathcal{L}(x, \lambda_1) = f(x) - \lambda_1 c_1(x)$$

At solution x^* , there is a scalar λ_1^* such that $\nabla_x \mathcal{L}(x^*, \lambda_1^*) = 0$

This condition is **necessary** but **not sufficient**.

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\mathcal{I} = \{1\}, \quad \mathcal{E} = \emptyset$$

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\mathcal{I} = \{1\}, \quad \mathcal{E} = \emptyset$$

Q: What is feasible set?



Constrained Optimization

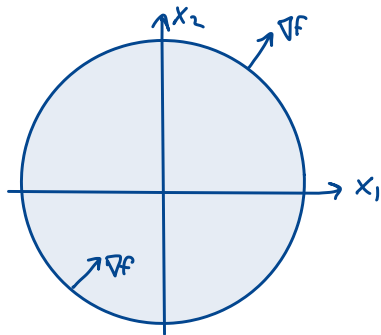
A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\mathcal{I} = \{1\}, \quad \mathcal{E} = \emptyset$$



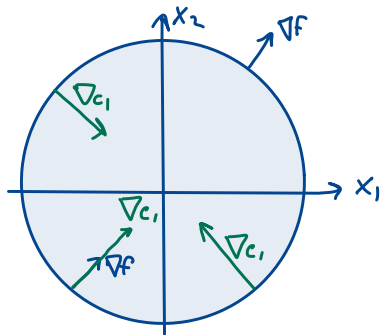
Q: What is feasible set?

A: Now, feasible set consists of the circle and its interior!

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$



$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\mathcal{I} = \{1\}, \quad \mathcal{E} = \emptyset$$

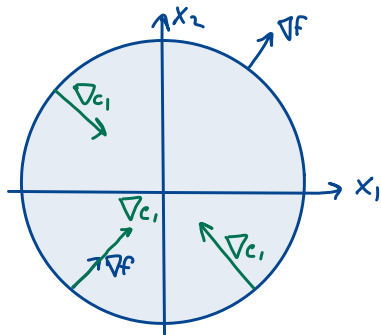
$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Constraint normal ∇c_1 points toward the interior of the feasible region at each point on the boundary of the circle.

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$



$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\mathcal{I} = \{1\}, \quad \mathcal{E} = \emptyset$$

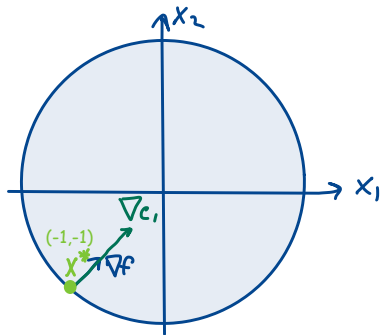
$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Q: What is the solution x^* ?

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$



$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\mathcal{I} = \{1\}, \quad \mathcal{E} = \emptyset$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Q: What is the solution x^* ?

$$\mathbf{A: } x^* = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

A given feasible point \mathbf{x} is **not optimal**, if we can find a small step \mathbf{s} that **both**

- retains feasibility,
- decreases the objective function $f(x)$ to first order.

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

A given feasible point \mathbf{x} is **not optimal**, if we can find a small step \mathbf{s} that **both**

- retains feasibility,
- decreases the objective function $f(x)$ to first order.

Approximate $c_1(x)$ to first order: $c_1(x + \mathbf{s}) \approx c_1(x) + \nabla c_1(x)^\top \mathbf{s}$

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

A given feasible point \mathbf{x} is **not optimal**, if we can find a small step \mathbf{s} that **both**

- retains feasibility,
- decreases the objective function $f(x)$ to first order.

Approximate $c_1(x)$ to first order: $c_1(x + \mathbf{s}) \approx c_1(x) + \nabla c_1(x)^\top \mathbf{s}$

If \mathbf{s} retains feasibility $\implies c_1(x) + \nabla c_1(x)^\top \mathbf{s} \geq 0$

Constrained Optimization

A Single Inequality Constraint

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

A given feasible point \mathbf{x} is **not optimal**, if we can find a small step \mathbf{s} that **both**

- retains feasibility, $\implies c_1(x) + \nabla c_1(x)^\top \mathbf{s} \geq 0$
- decreases the objective function $f(x)$ to first order.

Similarly, approximate $f(x)$ to first order: $f(x + \mathbf{s}) \approx f(x) + \nabla f(x)^\top \mathbf{s}$

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

A given feasible point \mathbf{x} is **not optimal**, if we can find a small step \mathbf{s} that **both**

- retains feasibility, $\implies c_1(x) + \nabla c_1(x)^\top \mathbf{s} \geq 0$
- decreases the objective function $f(x)$ to first order.

Similarly, approximate $f(x)$ to first order: $f(x + \mathbf{s}) \approx f(x) + \nabla f(x)^\top \mathbf{s}$

$f(x)$ is decreasing $\implies f(x + \mathbf{s}) - f(x) < 0$

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

A given feasible point \mathbf{x} is **not optimal**, if we can find a small step \mathbf{s} that **both**

- retains feasibility, $\implies c_1(x) + \nabla c_1(x)^\top \mathbf{s} \geq 0$
- decreases the objective function $f(x)$ to first order.

Similarly, approximate $f(x)$ to first order: $f(x + \mathbf{s}) \approx f(x) + \nabla f(x)^\top \mathbf{s}$

$f(x)$ is decreasing $\implies f(x + \mathbf{s}) - f(x) < 0$

$$f(x) + \nabla f(x)^\top \mathbf{s} - f(x) < 0$$

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

A given feasible point \mathbf{x} is **not optimal**, if we can find a small step \mathbf{s} that **both**

- retains feasibility, $\implies c_1(x) + \nabla c_1(x)^\top \mathbf{s} \geq 0$
- decreases the objective function $f(x)$ to first order.

Similarly, approximate $f(x)$ to first order: $f(x + \mathbf{s}) \approx f(x) + \nabla f(x)^\top \mathbf{s}$

$f(x)$ is decreasing $\implies f(x + \mathbf{s}) - f(x) < 0$

$$f(x) + \nabla f(x)^\top \mathbf{s} - f(x) < 0 \implies \nabla f(x)^\top \mathbf{s} < 0$$

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

A given feasible point \mathbf{x} is **not optimal**, if we can find a small step \mathbf{s} that **both**

C1: • retains feasibility, $\implies c_1(x) + \nabla c_1(x)^\top \mathbf{s} \geq 0$

C2: • decreases the objective function $f(x)$ to first order. $\implies \nabla f(x)^\top \mathbf{s} < 0$

Constrained Optimization

A Single Inequality Constraint

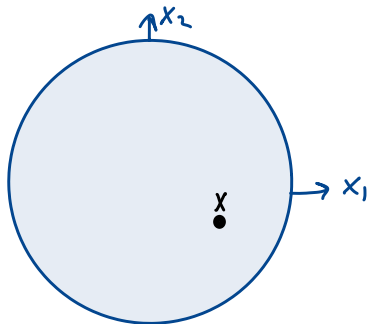
$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Case 1: Given x lies strictly inside the circle, $c_1(x) > 0$



Q: How would you select s ?

Remember the conditions:

$$\mathbf{C1:} \quad c_1(x) + \nabla c_1(x)^\top s \geq 0$$

$$\mathbf{C2:} \quad \nabla f(x)^\top s < 0$$

Constrained Optimization

A Single Inequality Constraint

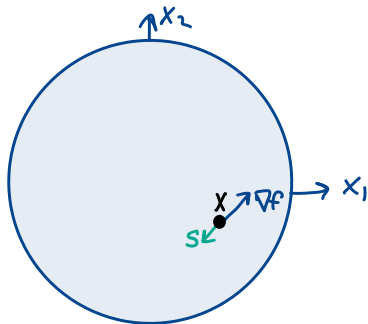
$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Case 1: Given x lies strictly inside the circle, $c_1(x) > 0$



Q: How would you select s ?

$$s = -\alpha \nabla f(x)$$

for any positive scalar α
sufficiently small.

Remember the conditions:

$$\text{C1: } c_1(x) + \nabla c_1(x)^\top s \geq 0$$

$$\text{C2: } \nabla f(x)^\top s < 0$$

Constrained Optimization

A Single Inequality Constraint

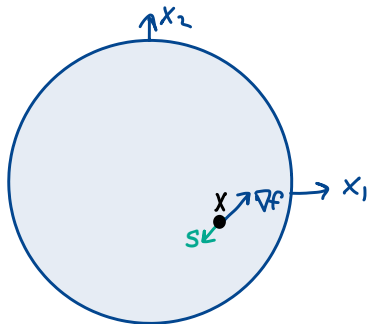
$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Case 1: Given x lies strictly inside the circle, $c_1(x) > 0$



Q: How would you select s ?

$$s = -\alpha \nabla f(x)$$

for any positive scalar α
sufficiently small.

However, no step s is given
when $\nabla f(x) = 0$

Remember the conditions:

$$\text{C1: } c_1(x) + \nabla c_1(x)^\top s \geq 0$$

$$\text{C2: } \nabla f(x)^\top s < 0$$

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Case 2: Given \mathbf{x} lies on the boundary of the circle, $c_1(\mathbf{x}) = 0$

Remember **C1**: $c_1(\mathbf{x}) + \nabla c_1(\mathbf{x})^\top \mathbf{s} \geq 0$.

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Case 2: Given x lies on the boundary of the circle, $c_1(x) = 0$

Remember **C1**: $c_1(x) + \nabla c_1(x)^T s \geq 0$.

C1: $\nabla c_1(x)^T s \geq 0$

C2: $\nabla f(x)^T s < 0$

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

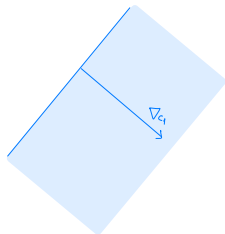
$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Case 2: Given x lies on the boundary of the circle, $c_1(x) = 0$

C1: $\nabla c_1(x)^\top s \geq 0 \rightarrow$ *Closed half-space*



Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

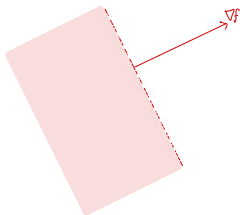
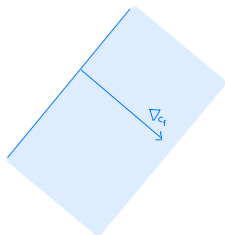
$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Case 2: Given \mathbf{x} lies on the boundary of the circle, $c_1(\mathbf{x}) = 0$

C1: $\nabla c_1(\mathbf{x})^\top \mathbf{s} \geq 0 \rightarrow$ *Closed half-space*

C2: $\nabla f(\mathbf{x})^\top \mathbf{s} < 0 \rightarrow$ *Open half-space*



Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

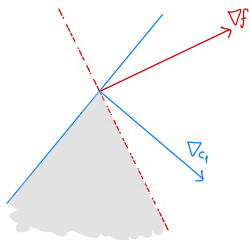
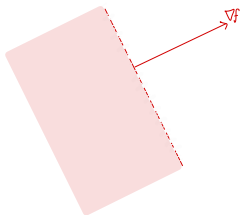
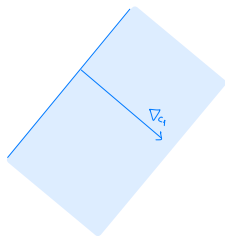
$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Case 2: Given \mathbf{x} lies on the boundary of the circle, $c_1(\mathbf{x}) = 0$

C1: $\nabla c_1(\mathbf{x})^\top \mathbf{s} \geq 0 \rightarrow$ *Closed half-space*

C2: $\nabla f(\mathbf{x})^\top \mathbf{s} < 0 \rightarrow$ *Open half-space*



Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Case 2: Given \mathbf{x} lies on the boundary of the circle, $c_1(\mathbf{x}) = 0$

If ∇f and ∇c_1 point in the opposite direction

$$\nabla f = \lambda_1 \nabla c_1 \text{ for some } \lambda_1 < 0$$

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

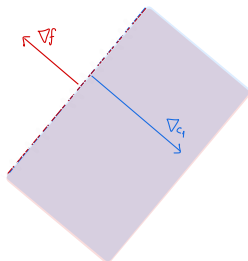
Case 2: Given \mathbf{x} lies on the boundary of the circle, $c_1(\mathbf{x}) = 0$

If ∇f and ∇c_1 point in the opposite direction

$$\nabla f = \lambda_1 \nabla c_1 \quad \text{for some } \lambda_1 < 0$$

Intersection region is

entire open half-space!



Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Case 2: Given \mathbf{x} lies on the boundary of the circle, $c_1(\mathbf{x}) = 0$

If ∇f and ∇c_1 point in the same direction

$$\nabla f = \lambda_1 \nabla c_1 \text{ for some } \lambda_1 \geq 0$$

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

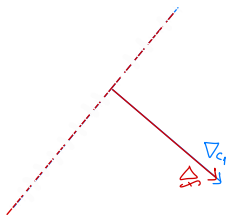
$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Case 2: Given x lies on the boundary of the circle, $c_1(x) = 0$

If ∇f and ∇c_1 point in the same direction

$$\nabla f = \lambda_1 \nabla c_1 \quad \text{for some } \lambda_1 \geq 0$$

Intersection region is **empty!**



Constrained Optimization

A Single Inequality Constraint

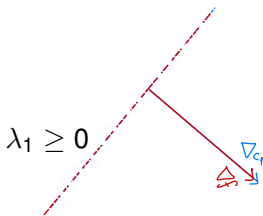
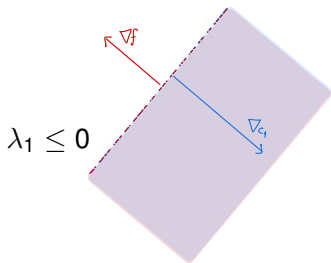
$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Case 2: Given x lies on the boundary of the circle, $c_1(x) = 0$



Q: Which one shows the convergence?

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Case 1: Given x lies strictly inside the circle, $c_1(x) > 0$

Case 2: Given x lies on the boundary of the circle, $c_1(x) = 0$

Optimality Conditions for both Case 1 and Case 2:

When no first order feasible descent direction exists at some point x^* , we have that

$$\nabla_x \mathcal{L}(x^*, \lambda_1^*) = 0 \text{ for some } \lambda_1^* \geq 0.$$

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Case 1: Given \mathbf{x} lies strictly inside the circle, $c_1(\mathbf{x}) > 0$

Case 2: Given \mathbf{x} lies on the boundary of the circle, $c_1(\mathbf{x}) = 0$

Optimality Conditions for both Case 1 and Case 2:

When no first order feasible descent direction exists at some point \mathbf{x}^* , we have that

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \lambda_1^*) = 0 \text{ for some } \lambda_1^* \geq 0.$$

We also require: $\lambda_1^* c_1(\mathbf{x}^*) = 0 \rightarrow$ *Complementarity Condition*

Constrained Optimization

A Single Inequality Constraint

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0$$

$$f(x) = x_1 + x_2$$

$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

Case 1: Given x lies strictly inside the circle, $c_1(x) > 0$

Case 2: Given x lies on the boundary of the circle, $c_1(x) = 0$

Optimality Conditions for both Case 1 and Case 2:

When no first order feasible descent direction exists at some point x^* , we have that

$$\nabla_x \mathcal{L}(x^*, \lambda_1^*) = 0 \text{ for some } \lambda_1^* \geq 0.$$

We also require: $\lambda_1^* c_1(x^*) = 0 \rightarrow$ *Complementarity Condition*

λ_1 can be strictly positive **only** when the corresponding c_1 is **active**.

Constrained Optimization

$$\begin{aligned}\nabla_x \mathcal{L}(x^*, \lambda^*) &= 0, \\ c_i(x^*) &= 0, \quad \text{for all } i \in \mathcal{E}, \\ c_i(x^*) &\geq 0, \quad \text{for all } i \in \mathcal{I}, \\ \lambda_i^* &\geq 0, \quad \text{for all } i \in \mathcal{I}, \\ \lambda_i^* c_i(x^*) &= 0, \quad \text{for all } i \in \mathcal{E} \cup \mathcal{I}.\end{aligned}$$

Constrained Optimization

$$\begin{aligned}\nabla_x \mathcal{L}(x^*, \lambda^*) &= 0, \\ c_i(x^*) &= 0, \quad \text{for all } i \in \mathcal{E}, \\ c_i(x^*) &\geq 0, \quad \text{for all } i \in \mathcal{I}, \\ \lambda_i^* &\geq 0, \quad \text{for all } i \in \mathcal{I}, \\ \lambda_i^* c_i(x^*) &= 0, \quad \text{for all } i \in \mathcal{E} \cup \mathcal{I}.\end{aligned}$$

Often known as the **Karush-Kuhn-Tucker (KKT)** conditions.

Summary

Summary

- ▶ Safety in RL is an *active* and *popular* research area.

Summary

- ▶ **Safety in RL** is an *active* and *popular* research area.
- ▶ **Definitions** and **methodologies** are subject to change depending on the applications and requirements.



Summary

- ▶ **Safety in RL** is an *active* and *popular* research area.
- ▶ **Definitions** and **methodologies** are subject to change depending on the applications and requirements.
- ▶ Adapting **optimization procedure** to safety requirements are often preferred, especially for a *known / partially known* transition dynamics and environment.



Summary

- ▶ **Safety in RL** is an *active* and *popular* research area.
- ▶ **Definitions** and **methodologies** are subject to change depending on the applications and requirements.
- ▶ Adapting **optimization procedure** to safety requirements are often preferred, especially for a *known / partially known* transition dynamics and environment.
- ▶ This adaptation for **constrained optimal control** should be performed in such a way that **the KKT conditions must be satisfied**.