# A!

**Aalto University
School of Electrical
Engineering**

# Safety and Constrained Optimal Control

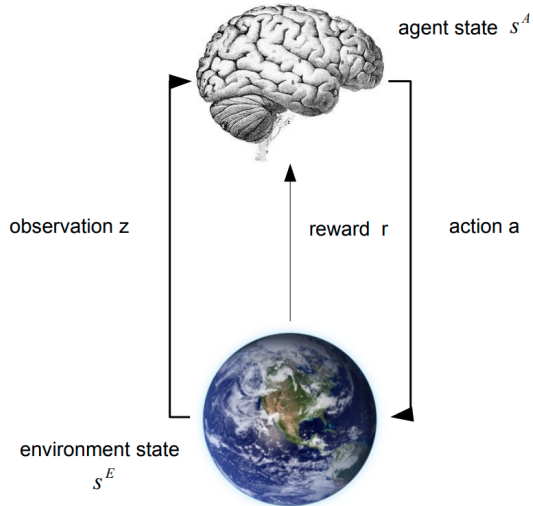Gökhan Alcan

📍 Dept. of Electrical Engineering and Automation
✉ gokhan.alcan@aalto.fi
🌐 www.gokhanalcan.com

November 3, 2020

# Reinforcement Learning



agent state $s^A$

observation z

reward r

action a

environment state

$s^E$

# 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!

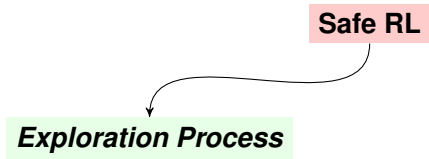► The agent is trained to *maximize the expected return* in a given task *while not taking any action* that *gives damage* to the environment or itself during learning and/or deployment.
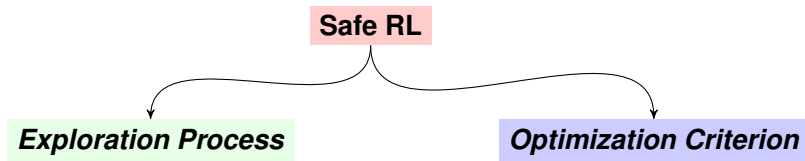
# Safety in Reinforcement Learning

Safe RL

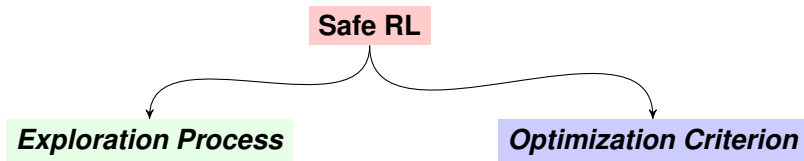# Safety in Reinforcement Learning

**Safe RL**

*Exploration Process*

# Safety in Reinforcement Learning

```
                    ┌──────────┐
                    │  Safe RL │
                    └──────────┘
              ┌───────────┴───────────┐
              ▼                       ▼
    ┌────────────────────┐  ┌────────────────────────┐
    │ Exploration Process│  │ Optimization Criterion │
    └────────────────────┘  └────────────────────────┘
```

# Safety in Reinforcement Learning

**Safe RL**

*Exploration Process*

*Optimization Criterion*

▶ Risk-directed Exploration

# Safety in Reinforcement Learning

**Safe RL**

**Exploration Process**

**Optimization Criterion**

▶ Risk-directed Exploration

▶ Utilization of External Knowledge

# Safety in Reinforcement Learning

**Safe RL**

**Exploration Process**

► Risk-directed Exploration

► Utilization of External Knowledge

**Optimization Criterion**

► Constrained Criterion

# Safety in Reinforcement Learning

**Safe RL**

**_Exploration Process_**

- Risk-directed Exploration
- Utilization of External Knowledge

**_Optimization Criterion_**

- Constrained Criterion
- Worst Case Criterion

# Safety in Reinforcement Learning

**Safe RL**

**Exploration Process**

▶ Risk-directed Exploration

▶ Utilization of External Knowledge

**Optimization Criterion**

▶ Constrained Criterion

▶ Worst Case Criterion

▶ Risk-Sensitive Criterion

# Safe Exploration

*OpenAI Safety-Gym*



**OpenAI Safety-Gym:** A. Ray, J. Achiam, and D. Amodei, "Benchmarking Safe Exploration in Deep Reinforcement Learning," 2019, https://cdn.openai.com/safexp-short.pdf.
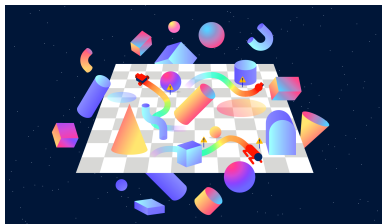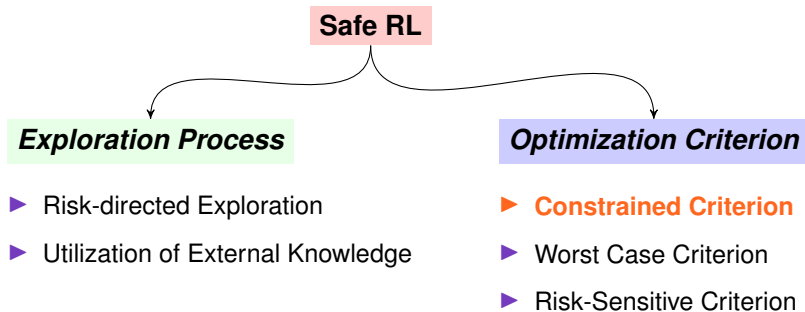
# Safe Exploration

## *OpenAI Safety-Gym*



## *Some Methods*

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

# Safety in Reinforcement Learning

**Safe RL**

**Exploration Process**

▶ Risk-directed Exploration

▶ Utilization of External Knowledge

**Optimization Criterion**

▶ **Constrained Criterion**

▶ Worst Case Criterion

▶ Risk-Sensitive Criterion

# Constrained Optimal Control

# Constrained Optimization

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

# Constrained Optimization

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

**Feasible Set:**

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

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

# Constrained Optimization

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

**Feasible Set:**

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

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

**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*

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

$$f(x) = x_1 + x_2$$
$$c_1(x) = x_1^2 + x_2^2 - 2$$
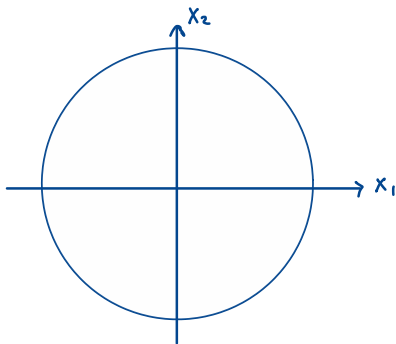$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

# 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*

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



$$f(x) = x_1 + x_2$$
$$c_1(x) = x_1^2 + x_2^2 - 2$$
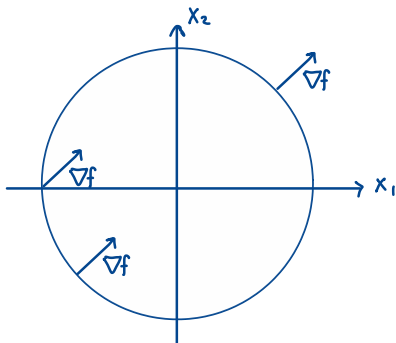$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

**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*

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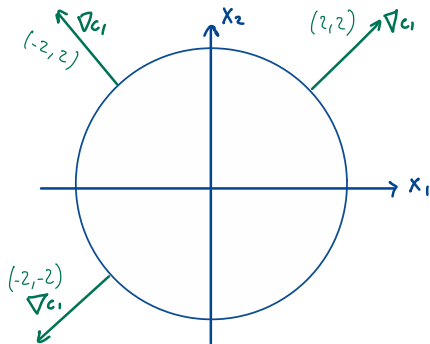


$$f(x) = x_1 + x_2$$
$$c_1(x) = x_1^2 + x_2^2 - 2$$
$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

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

# Constrained Optimization

*A Single Equality Constraint*

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



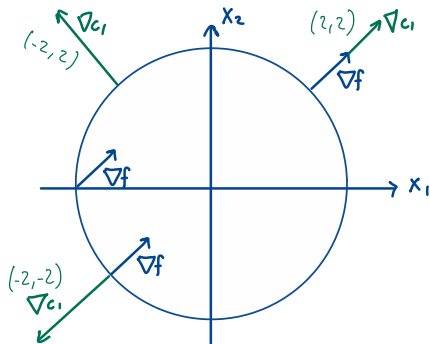$$f(x) = x_1 + x_2$$
$$c_1(x) = x_1^2 + x_2^2 - 2$$
$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

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

# Constrained Optimization

*A Single Equality Constraint*

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



$$f(x) = x_1 + x_2$$
$$c_1(x) = x_1^2 + x_2^2 - 2$$
$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$

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

**Q: What is the solution $x^*$?**

# Constrained Optimization

*A Single Equality Constraint*

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

$$f(x) = x_1 + x_2$$
$$c_1(x) = x_1^2 + x_2^2 - 2$$
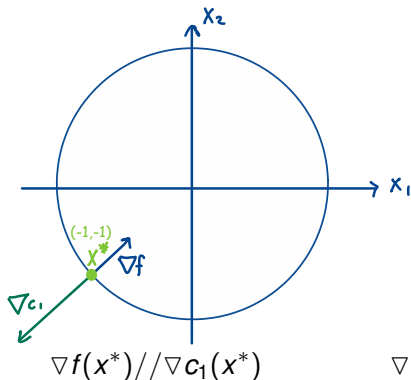$$\mathcal{I} = \emptyset, \quad \mathcal{E} = \{1\}$$



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

**Q: What is the solution $x^*$?**

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

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

# 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 $\lambda_1^*$ 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 $\lambda_1^*$ such that $\nabla_x \mathcal{L}(x^*, \lambda_1^*) = 0$

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

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

# 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 $\lambda_1^*$ such that $\nabla_x \mathcal{L}(x^*, \lambda_1^*) = 0$

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

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

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 $\lambda_1^*$ such that $\nabla_x \mathcal{L}(x^*, \lambda_1^*) = 0$

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

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

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 $\lambda_1^*$ such that $\nabla_x \mathcal{L}(x^*, \lambda_1^*) = 0$

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

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

**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 $\lambda_1^*$ such that $\nabla_x \mathcal{L}(x^*, \lambda_1^*) = 0$

This condition is **necessary** but **not sufficient**.
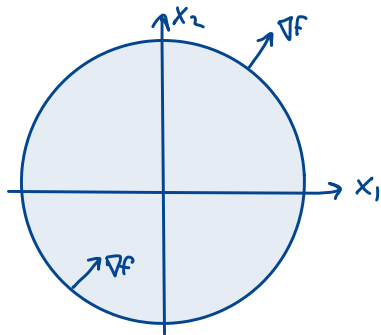
# Constrained Optimization

## *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*

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

$$f(x) = x_1 + x_2$$
$$c_1(x) = 2 - x_1^2 - x_2^2$$
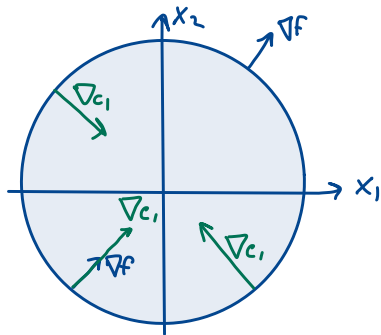$$\mathcal{I} = \{1\}, \quad \mathcal{E} = \emptyset$$

**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

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



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

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

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

$$\nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \qquad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$
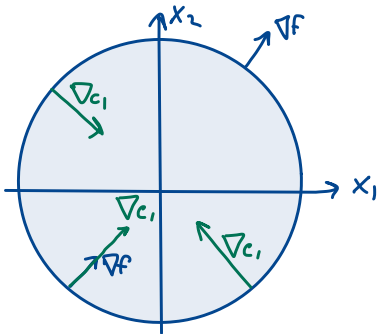
*Constraint normal $\nabla c_1$ points toward the interior of the feasible region at each point on the boundary of the circle.*

# Constrained Optimization

## *A Single Inequality Constraint*

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



$$f(x) = x_1 + x_2$$
$$c_1(x) = 2 - x_1^2 - x_2^2$$
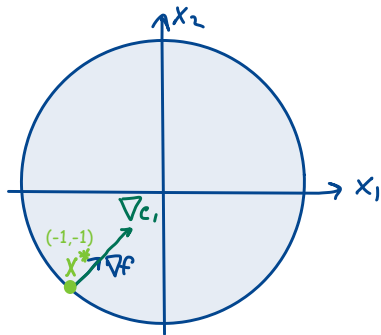$$\mathcal{I} = \{1\}, \quad \mathcal{E} = \emptyset$$

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

**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 $x^*$?**

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

# Constrained Optimization

*A Single Inequality Constraint*

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

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

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

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

# Constrained Optimization

## *A Single Inequality Constraint*

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

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

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

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

- retains feasibility,

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

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

# Constrained Optimization

*A Single Inequality Constraint*

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

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

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

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

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

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

# Constrained Optimization

### *A Single Inequality Constraint*

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

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

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

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

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

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

# Constrained Optimization

*A Single Inequality Constraint*

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

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

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

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

- retains feasibility, $\implies c_1(x) + \nabla c_1(x)^\top s \geq 0$

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

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

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

# Constrained Optimization

### *A Single Inequality Constraint*

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

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

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

- retains feasibility, $\implies c_1(x) + \nabla c_1(x)^\top s \geq 0$

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

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

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

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

# Constrained Optimization

## *A Single Inequality Constraint*

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

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

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

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

- retains feasibility, $\implies c_1(x) + \nabla c_1(x)^\top s \geq 0$

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

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

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

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

# Constrained Optimization

*A Single Inequality Constraint*

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

$$\min_{x_1,x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

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

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

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

# Constrained Optimization

## *A Single Inequality Constraint*

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

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

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \qquad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

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



**Q: How would you select s?**

*Remember the conditions:*
**C1:** $c_1(x) + \nabla c_1(x)^\top s \geq 0$
**C2:** $\nabla f(x)^\top s < 0$

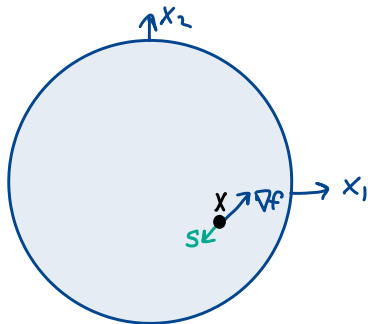# Constrained Optimization
## *A Single Inequality Constraint*

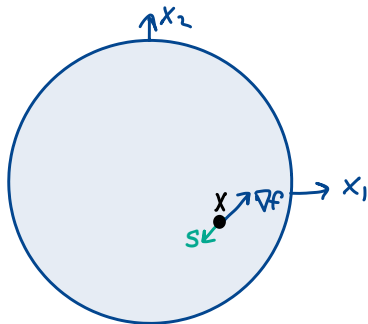$$f(x) = x_1 + x_2$$
$$c_1(x) = 2 - x_1^2 - x_2^2$$

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

**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 $\alpha$
sufficiently small.

*Remember the conditions:*
**C1:** $c_1(x) + \nabla c_1(x)^\top s \geq 0$
**C2:** $\nabla f(x)^\top s < 0$

# Constrained Optimization

## *A Single Inequality Constraint*

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

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

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \qquad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

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



**Q: How would you select s?**

$$\mathbf{s} = -\alpha \nabla f(x)$$
for any positive scalar $\alpha$ sufficiently small.

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

*Remember the conditions:*
**C1:** $c_1(x) + \nabla c_1(x)^\top s \geq 0$
**C2:** $\nabla f(x)^\top s < 0$

# Constrained Optimization

## *A Single Inequality Constraint*

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

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

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

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

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

# Constrained Optimization

### *A Single Inequality Constraint*

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

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

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

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

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

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

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

# Constrained Optimization

### *A Single Inequality Constraint*

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

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

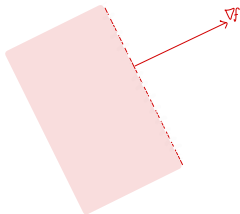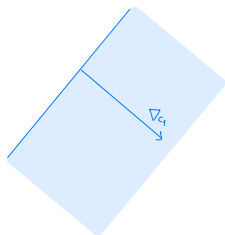$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

**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*

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

$$\min_{x_1,x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

**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*

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

# Constrained Optimization

**A Single Inequality Constraint**

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

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

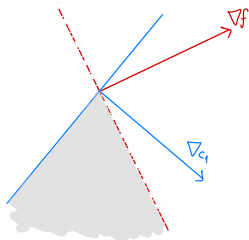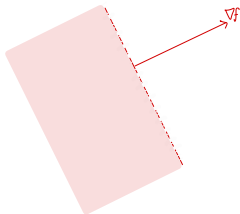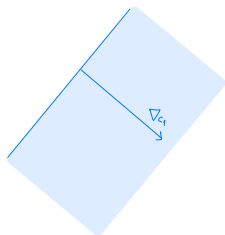$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

**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*

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

# Constrained Optimization

### *A Single Inequality Constraint*

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

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

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

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

If $\nabla f$ and $\nabla c_1$ point in the opposite direction

$\nabla f = \lambda_1 \nabla c_1$ for some $\lambda_1 < 0$

# Constrained Optimization

**A Single Inequality Constraint**

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

$$\min_{x_1,x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$
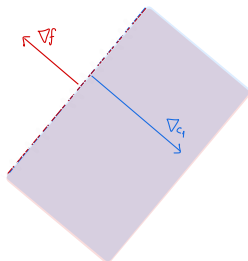
**Case 2:** Given **x** lies on the boundary of the circle, $c_1(x) = 0$

If $\nabla f$ and $\nabla c_1$ point in the opposite direction

$\nabla f = \lambda_1 \nabla c_1$ for some $\lambda_1 < 0$

Intersection region is
**entire open half-space!**

# Constrained Optimization

### *A Single Inequality Constraint*

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

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

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

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

If $\nabla f$ and $\nabla c_1$ point in the same direction

$\nabla f = \lambda_1 \nabla c_1$ for some $\lambda_1 \geq 0$

# Constrained Optimization

## *A Single Inequality Constraint*

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

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

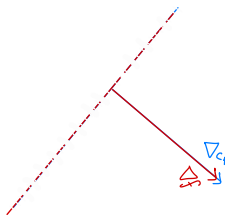$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

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

If $\nabla f$ and $\nabla c_1$ point in the same direction

$\nabla f = \lambda_1 \nabla c_1$ for some $\lambda_1 \geq 0$

Intersection region is **empty!**

# Constrained Optimization

*A Single Inequality Constraint*

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

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

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

**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

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

# Constrained Optimization

### *A Single Inequality Constraint*

$f(x) = x_1 + x_2$

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

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

**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

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

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

# Constrained Optimization

*A Single Inequality Constraint*

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

$$\min_{x_1, x_2} x_1 + x_2 \quad \text{s.t.} \quad 2 - x_1^2 - x_2^2 \geq 0 \qquad \nabla f = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \nabla c_1 = \begin{bmatrix} -2x_1 \\ -2x_2 \end{bmatrix}$$

**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*

$\lambda_1$ can be strictly positive **only** when the corresponding $c_1$ is **active**.

# Constrained Optimization

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

# Constrained Optimization

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

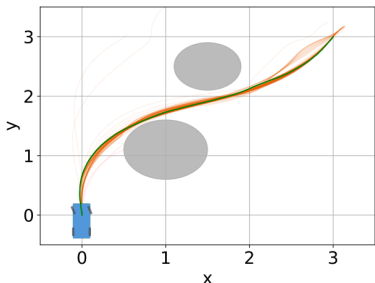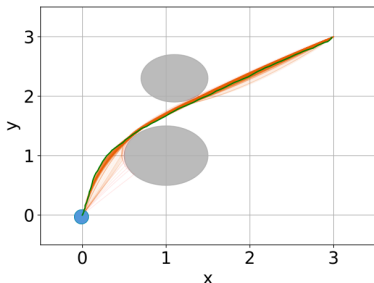Often known as the Karush-Kuhn-Tucker (**KKT**) conditions.

# Constrained Optimization

*Robotic Application:* *Safe Trajectory Optimization*

$$\min_{\mathbf{u}_0,\ldots,\mathbf{u}_{N-1}} \quad \ell_f(\mathbf{x}_N) + \sum_{k=0}^{N-1} \ell(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{subject to} \quad \mathbf{x}_{k+1} = f(\mathbf{x}_k, \mathbf{u}_k),$$

$$\mathbf{g}(\mathbf{x}_k, \mathbf{u}_k) \geq \mathbf{0},$$

# Constrained Optimization

*Robotic Application:* *Safe Trajectory Optimization*

$$\min_{\mathbf{u}_0, \ldots, \mathbf{u}_{N-1}} \quad \ell_f(\mathbf{x}_N) + \sum_{k=0}^{N-1} \ell(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{subject to} \quad \mathbf{x}_{k+1} = f(\mathbf{x}_k, \mathbf{u}_k),$$

$$\mathbf{g}(\mathbf{x}_k, \mathbf{u}_k) \geq \mathbf{0},$$

# Summary

# Summary

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

**Aalto University**
**School of Electrical**
**Engineering**

# 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**.