



ΠΑΝΕΠΙΣΤΗΜΙΟ
ΠΑΤΡΩΝ
UNIVERSITY OF PATRAS

Optimization for Robotics Summer School

Lecture 1 & 2: Introduction to Optimization

Konstantinos Chatzilygeroudis - costashatz@upatras.gr

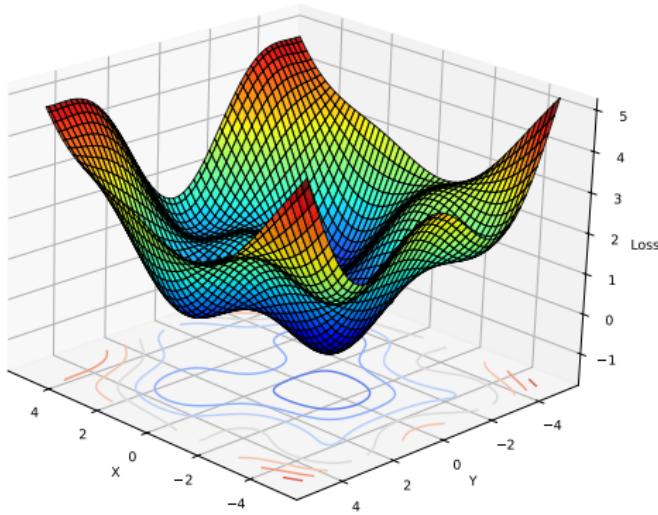
Department of Electrical and Computer Engineering
University of Patras

Template made by Panagiotis Papagiannopoulos

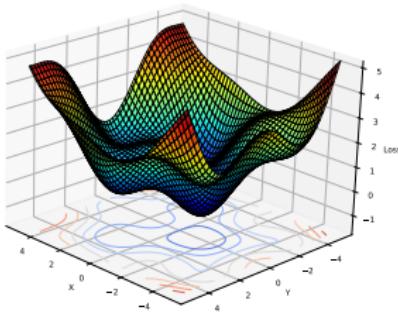


What is Optimization?

Optimization is the process of making something as effective or functional as possible.



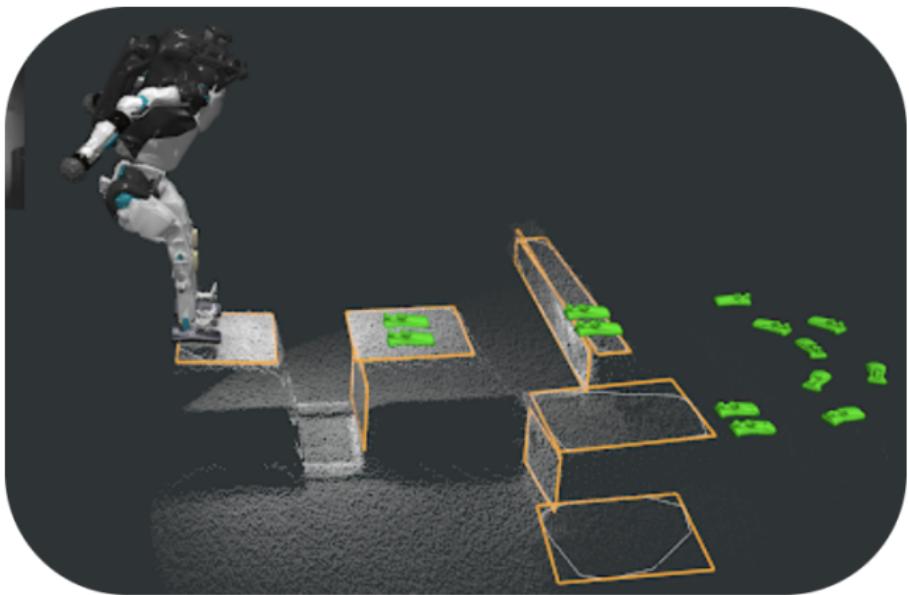
What is Optimization?



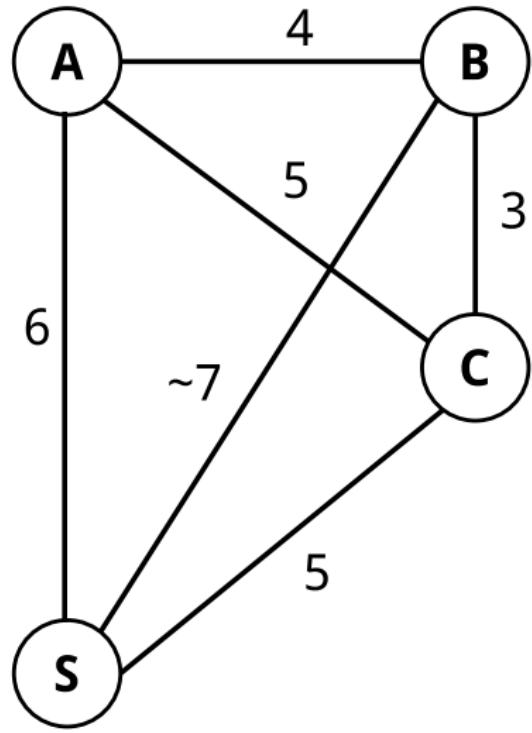
What is the goal of optimization?

- To find the best solution to a problem, which involves minimizing or maximizing a (set of) function(s) (**the objective function(s)**).
- To find the values of the *unknowns* (**the variables**) that lead to the optimal outcome, subject to certain **constraints**.

Why?

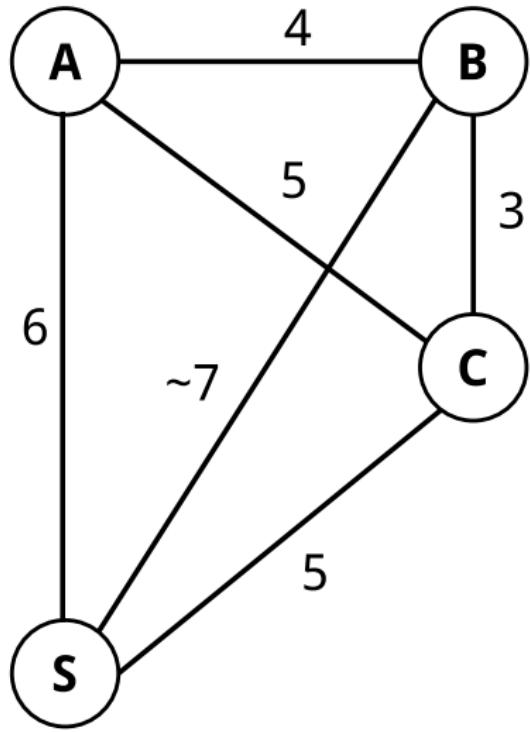


Discrete Optimization Example



Example: A delivery driver needs to drop off packages at 3 different houses, and they want to minimize total travel distance.

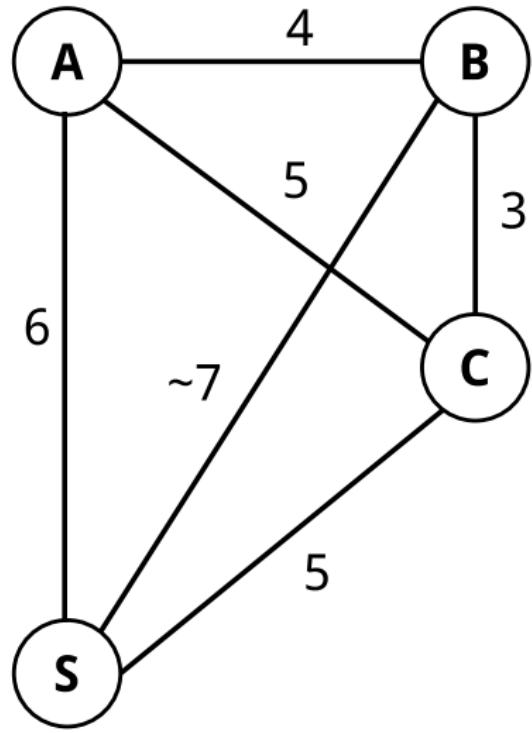
Discrete Optimization Example



Example: A delivery driver needs to drop off packages at 3 different houses, and they want to minimize total travel distance.

- The objective function is the total travel distance.
- We have a set of possible solutions. Can we enumerate them?

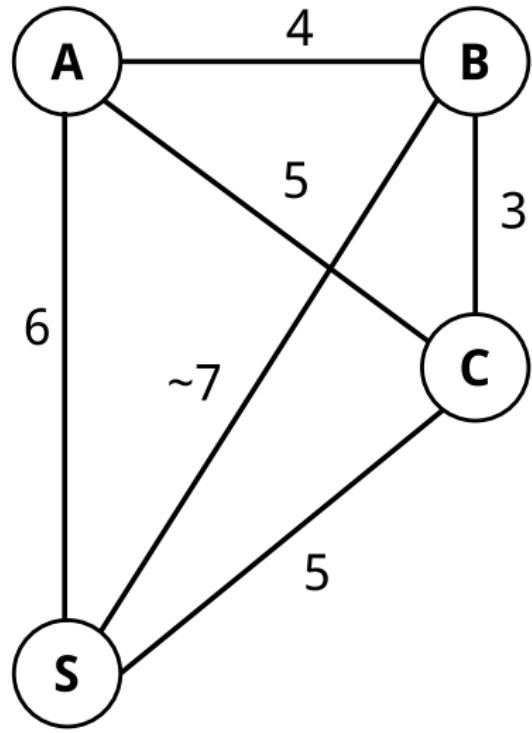
Discrete Optimization Example



Example: A delivery driver needs to drop off packages at 3 different houses, and they want to minimize total travel distance.

- The objective function is the total travel distance.
- We have a set of possible solutions. Can we enumerate them? **Yes!**

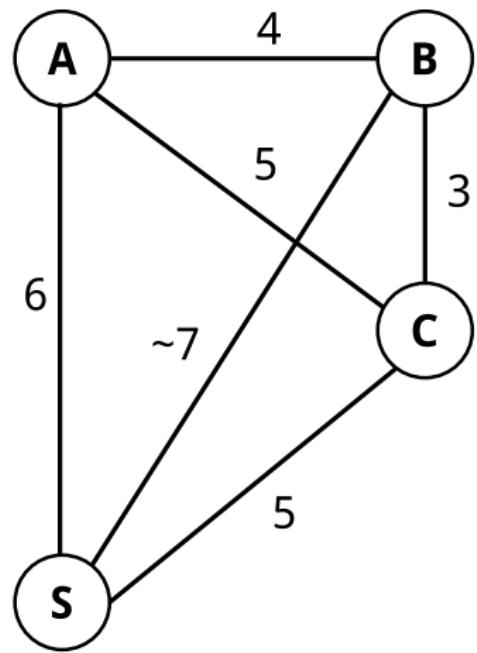
Discrete Optimization Example



Example: A delivery driver needs to drop off packages at 3 different houses, and they want to minimize total travel distance.

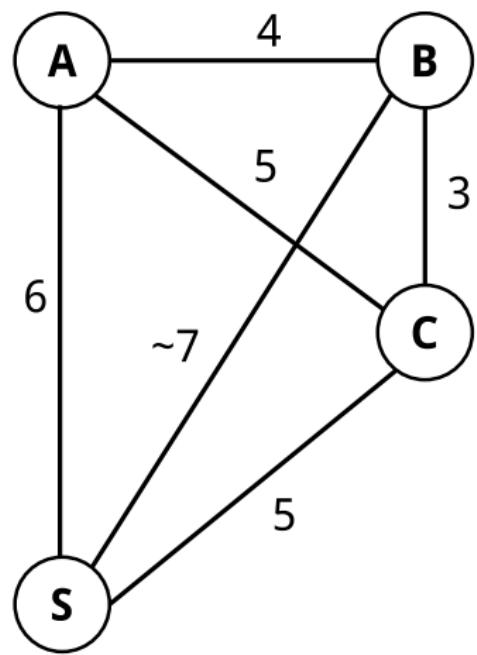
- The objective function is the total travel distance.
- We have a set of possible solutions. Can we enumerate them? **Yes!**
- What is the optimal solution?

Discrete Optimization Example (2)



Let's enumerate all possible solutions:

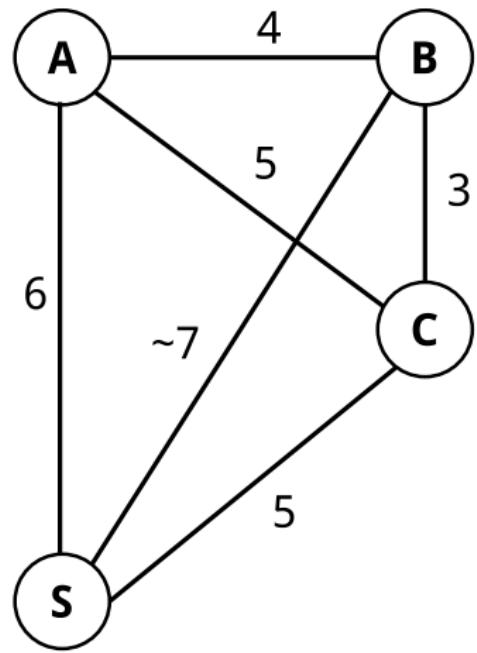
Discrete Optimization Example (2)



Let's enumerate all possible solutions:

- 1 $S \rightarrow A \rightarrow B \rightarrow C \rightarrow S$ (distance: 18)

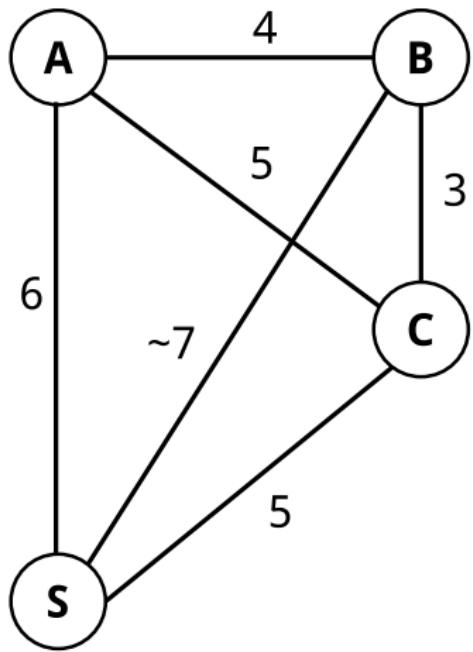
Discrete Optimization Example (2)



Let's enumerate all possible solutions:

- 1 $S \rightarrow A \rightarrow B \rightarrow C \rightarrow S$ (distance: 18)
- 2 $S \rightarrow A \rightarrow C \rightarrow B \rightarrow S$ (distance: 21)

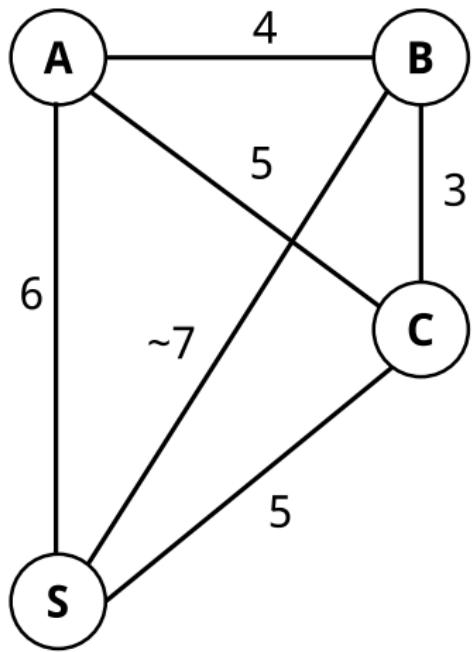
Discrete Optimization Example (2)



Let's enumerate all possible solutions:

- 1 $S \rightarrow A \rightarrow B \rightarrow C \rightarrow S$ (distance: 18)
- 2 $S \rightarrow A \rightarrow C \rightarrow B \rightarrow S$ (distance: 21)
- 3 $S \rightarrow B \rightarrow A \rightarrow C \rightarrow S$ (distance: 21)
- 4 $S \rightarrow B \rightarrow C \rightarrow A \rightarrow S$ (distance: 21)
- 5 $S \rightarrow C \rightarrow A \rightarrow B \rightarrow S$ (distance: 21)
- 6 $S \rightarrow C \rightarrow B \rightarrow A \rightarrow S$ (distance: 18)

Discrete Optimization Example (2)

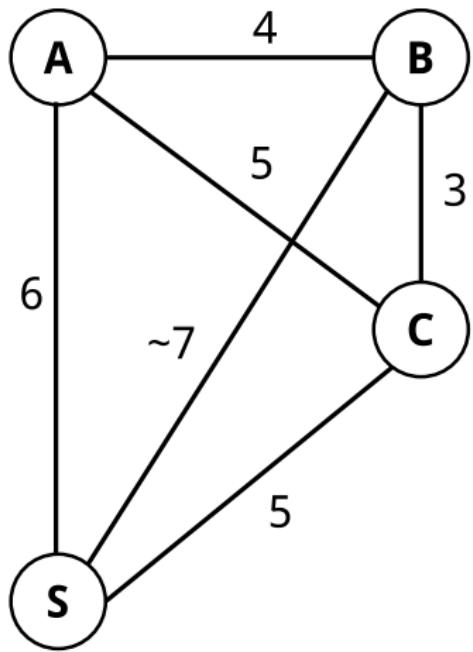


Let's enumerate all possible solutions:

- 1 $S \rightarrow A \rightarrow B \rightarrow C \rightarrow S$ (distance: 18)
- 2 $S \rightarrow A \rightarrow C \rightarrow B \rightarrow S$ (distance: 21)
- 3 $S \rightarrow B \rightarrow A \rightarrow C \rightarrow S$ (distance: 21)
- 4 $S \rightarrow B \rightarrow C \rightarrow A \rightarrow S$ (distance: 21)
- 5 $S \rightarrow C \rightarrow A \rightarrow B \rightarrow S$ (distance: 21)
- 6 $S \rightarrow C \rightarrow B \rightarrow A \rightarrow S$ (distance: 18)

The optimal solution is?

Discrete Optimization Example (2)

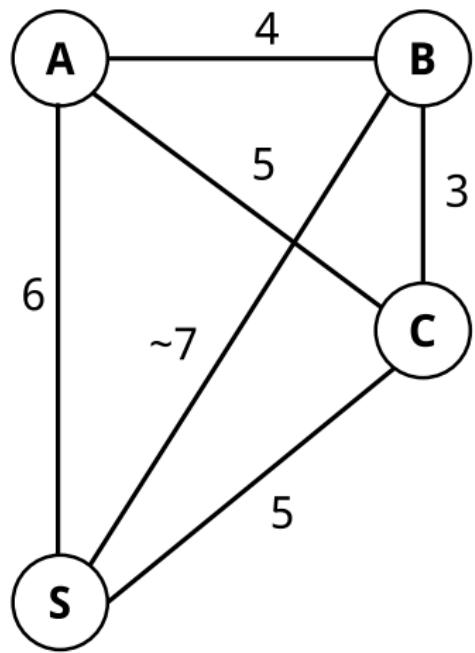


Let's enumerate all possible solutions:

- 1 $S \rightarrow A \rightarrow B \rightarrow C \rightarrow S$ (distance: 18)
- 2 $S \rightarrow A \rightarrow C \rightarrow B \rightarrow S$ (distance: 21)
- 3 $S \rightarrow B \rightarrow A \rightarrow C \rightarrow S$ (distance: 21)
- 4 $S \rightarrow B \rightarrow C \rightarrow A \rightarrow S$ (distance: 21)
- 5 $S \rightarrow C \rightarrow A \rightarrow B \rightarrow S$ (distance: 21)
- 6 $S \rightarrow C \rightarrow B \rightarrow A \rightarrow S$ (distance: 18)

The optimal solution is? We have two!

Discrete Optimization Example (2)



Let's enumerate all possible solutions:

- 1 $S \rightarrow A \rightarrow B \rightarrow C \rightarrow S$ (distance: 18)
- 2 $S \rightarrow A \rightarrow C \rightarrow B \rightarrow S$ (distance: 21)
- 3 $S \rightarrow B \rightarrow A \rightarrow C \rightarrow S$ (distance: 21)
- 4 $S \rightarrow B \rightarrow C \rightarrow A \rightarrow S$ (distance: 21)
- 5 $S \rightarrow C \rightarrow A \rightarrow B \rightarrow S$ (distance: 21)
- 6 $S \rightarrow C \rightarrow B \rightarrow A \rightarrow S$ (distance: 18)

The optimal solution is? We have two!

In general, we cannot enumerate all possible solutions! What can we do?

Some notation

$$\mathbf{y} = f(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}^M,$$

$$\mathbf{x} \in \mathbb{R}^N, \mathbf{y} \in \mathbb{R}^M.$$

What is $\frac{\partial f}{\partial \mathbf{x}} = \frac{\partial \mathbf{y}}{\partial \mathbf{x}}$?

Some notation

$$\mathbf{y} = f(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}^M,$$

$$\mathbf{x} \in \mathbb{R}^N, \mathbf{y} \in \mathbb{R}^M.$$

What is $\frac{\partial f}{\partial \mathbf{x}} = \frac{\partial \mathbf{y}}{\partial \mathbf{x}}$?

$$\frac{\partial f}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{y}_1}{\partial \mathbf{x}_1} & \frac{\partial \mathbf{y}_1}{\partial \mathbf{x}_2} & \cdots & \frac{\partial \mathbf{y}_1}{\partial \mathbf{x}_N} \\ \frac{\partial \mathbf{y}_2}{\partial \mathbf{x}_1} & \frac{\partial \mathbf{y}_2}{\partial \mathbf{x}_2} & \cdots & \frac{\partial \mathbf{y}_2}{\partial \mathbf{x}_N} \\ \ddots & \ddots & \ddots & \ddots \\ \frac{\partial \mathbf{y}_M}{\partial \mathbf{x}_1} & \frac{\partial \mathbf{y}_M}{\partial \mathbf{x}_2} & \cdots & \frac{\partial \mathbf{y}_M}{\partial \mathbf{x}_N} \end{bmatrix} \in \mathbb{R}^{M \times N}$$

This is called the **Jacobian**.

Notation Continued

- When $M = 1$, then $\frac{\partial f}{\partial \mathbf{x}} \in \mathbb{R}^{1 \times N}$ (row vector),
- When $M = 1$, we usually denote the **gradient** of f with $\nabla_{\mathbf{x}} f$; in other words, $\nabla_{\mathbf{x}} f = \frac{\partial f}{\partial \mathbf{x}}^T \in \mathbb{R}^{N \times 1}$ (column vector),
- When $M = 1$, the **Hessian** of f is:

$$\nabla_{\mathbf{x}\mathbf{x}}^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_N} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_N} \\ \ddots & \ddots & \ddots & \ddots \\ \frac{\partial^2 f}{\partial x_N \partial x_1} & \frac{\partial^2 f}{\partial x_N \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_N^2} \end{bmatrix}$$

We want to compute the Taylor expansion of f around $\bar{\mathbf{x}}$:

$$f(\mathbf{x}) \approx f(\bar{\mathbf{x}}) + \frac{\partial f}{\partial \mathbf{x}} \Big|_{\bar{\mathbf{x}}} (\mathbf{x} - \bar{\mathbf{x}}) + \dots$$

We can write: $\mathbf{x} = \bar{\mathbf{x}} + \Delta \mathbf{x}$, **and thus:**

$$f(\mathbf{x}) = f(\bar{\mathbf{x}} + \Delta \mathbf{x}) \approx f(\bar{\mathbf{x}}) + \frac{\partial f}{\partial \mathbf{x}} \Big|_{\bar{\mathbf{x}}} \Delta \mathbf{x} + \dots$$

Unconstrained Optimization (minimization!):

$$\min_{\mathbf{x}} f(\mathbf{x})$$

- $\mathbf{x} \in \mathbb{R}^N$

Unconstrained Optimization (minimization!):

$$\min_{\mathbf{x}} f(\mathbf{x})$$

- $\mathbf{x} \in \mathbb{R}^N$: **variables**
- $f(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}$

Unconstrained Optimization (minimization!):

$$\min_{\mathbf{x}} f(\mathbf{x})$$

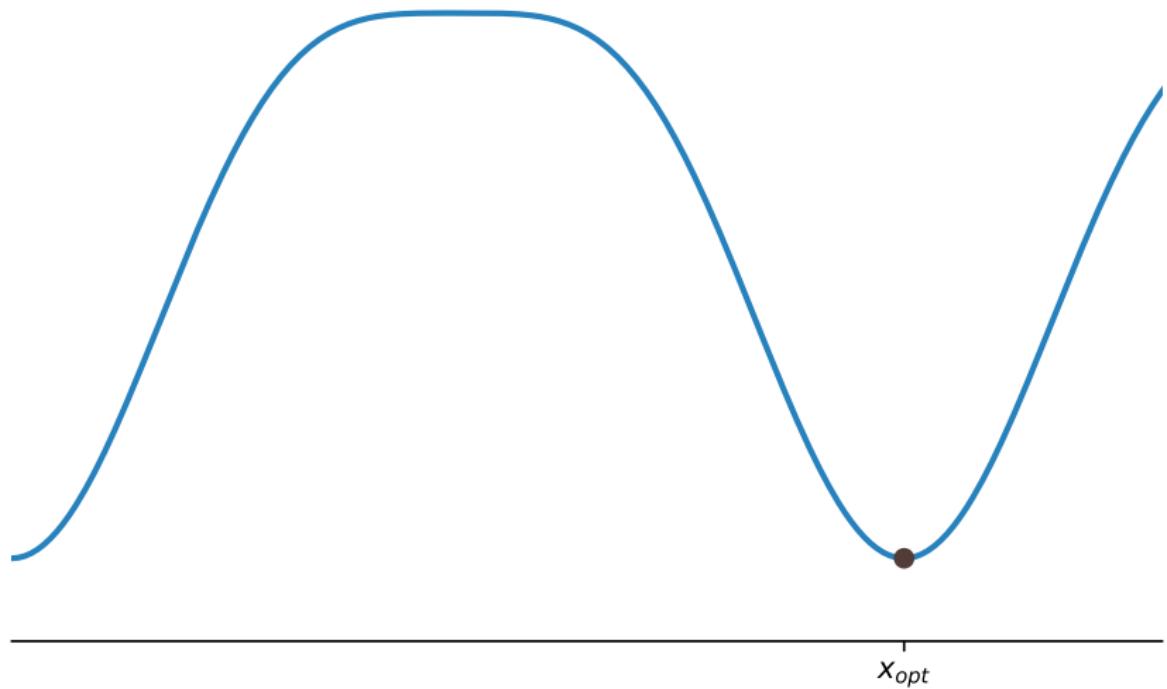
- $\mathbf{x} \in \mathbb{R}^N$: **variables**
- $f(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}$: **objective function**
- \mathbf{x}_* is the (global) **optimal solution**

Unconstrained Optimization (minimization!):

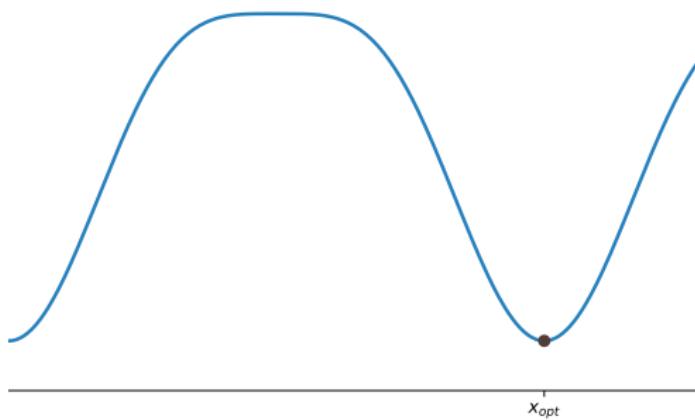
$$\min_{\mathbf{x}} f(\mathbf{x})$$

- $\mathbf{x} \in \mathbb{R}^N$: **variables**
- $f(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}$: **objective function**
- \mathbf{x}_* is the (global) **optimal solution**, if $f(\mathbf{x}_*) \leq f(\mathbf{x}), \forall \mathbf{x}$
- \mathbf{x}_* is a **local optimum**, if $f(\mathbf{x}_*) \leq f(\mathbf{x}), \forall \mathbf{x} \in \mathcal{N}_{\mathbf{x}_*}$ (neighborhood of \mathbf{x}_*)

How does a solution look like?



How does a solution look like? (2)



Necessary Condition for a (local) solution (when f is smooth):

$$\frac{\partial f}{\partial \mathbf{x}} \Big|_{\mathbf{x}_*} = 0 \Leftrightarrow \nabla_{\mathbf{x}} f(\mathbf{x}_*) = 0$$

How does a solution look like? (3)

First Order Necessary Condition for a (local) minimum:

$$\frac{\partial f}{\partial \mathbf{x}} \Big|_{\mathbf{x}_*} = 0 \Leftrightarrow \nabla_{\mathbf{x}} f(\mathbf{x}_*) = 0$$

How does a solution look like? (3)

First Order Necessary Condition for a (local) minimum:

$$\frac{\partial f}{\partial \mathbf{x}} \Big|_{\mathbf{x}_*} = 0 \Leftrightarrow \nabla_{\mathbf{x}} f(\mathbf{x}_*) = 0$$

Second Order Necessary Condition for a (local) minimum:

$$\nabla_{\mathbf{x}} f(\mathbf{x}_*) = 0$$

$$\nabla_{\mathbf{xx}}^2 f(\mathbf{x}_*) \text{ is psd}$$

Necessary and Sufficient Conditions

- **Necessary conditions** → if x_* is a (local) minimum, then these conditions hold;

Necessary and Sufficient Conditions

- **Necessary conditions** → if x_* is a (local) minimum, then these conditions hold;
- **Sufficient conditions**

- **Necessary conditions** → if x_* is a (local) minimum, then these conditions hold;
- **Sufficient conditions** → if these conditions hold for x_\dagger , then x_\dagger is a (local) minimum.

Enough of this! Let's see some action!

Problems related to optimization or maybe **useful to optimization**

Enough of this! Let's see some action!

Problems related to optimization or maybe **useful to optimization**

- **Root Finding:** Find \mathbf{x}^* such that $f(\mathbf{x}^*) = 0$,

Enough of this! Let's see some action!

Problems related to optimization or maybe useful to optimization

- **Root Finding:** Find \mathbf{x}^* such that $f(\mathbf{x}^*) = 0$,
- **Fixed Point:** Find \mathbf{x}^* such that $f(\mathbf{x}^*) = \mathbf{x}^*$,

Enough of this! Let's see some action!

Problems related to optimization or maybe useful to optimization

- **Root Finding:** Find \mathbf{x}^* such that $f(\mathbf{x}^*) = 0$,
- **Fixed Point:** Find \mathbf{x}^* such that $f(\mathbf{x}^*) = \mathbf{x}^*$,
- Many ways of solving this:
 - Fixed-Point Iteration
 - ...
 - Newton's Method

Newton's Method

- Linearize around current estimate, \mathbf{x}_k :

$$f(\mathbf{x}_k + \Delta\mathbf{x}) \approx f(\mathbf{x}_k) + \frac{\partial f}{\partial \mathbf{x}} \Big|_{\mathbf{x}_k} \Delta\mathbf{x}$$

- We then set $f(\mathbf{x}_k + \Delta\mathbf{x}) = 0$ and solve for $\Delta\mathbf{x}$:

$$\begin{aligned} f(\mathbf{x}_k) + \frac{\partial f}{\partial \mathbf{x}} \Big|_{\mathbf{x}_k} \Delta\mathbf{x} &= 0 \\ \Delta\mathbf{x} &= -\left(\frac{\partial f}{\partial \mathbf{x}} \Big|_{\mathbf{x}_k}\right)^{-1} f(\mathbf{x}_k) \end{aligned}$$

- We apply the correction $\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta\mathbf{x}$
- Repeat until convergence

And so?

$$\min_{\mathbf{x}} f(\mathbf{x}), f(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}$$

Necessary Condition for solution (when f is smooth):

$$\frac{\partial f}{\partial \mathbf{x}} \Big|_{\mathbf{x}^*} = 0$$

And so?

$$\min_{\mathbf{x}} f(\mathbf{x}), f(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}$$

Necessary Condition for solution (when f is smooth):

$$\frac{\partial f}{\partial \mathbf{x}} \Big|_{\mathbf{x}^*} = 0$$

Idea: We can frame optimization as root finding in the gradient!

Newton's Method for Optimization:

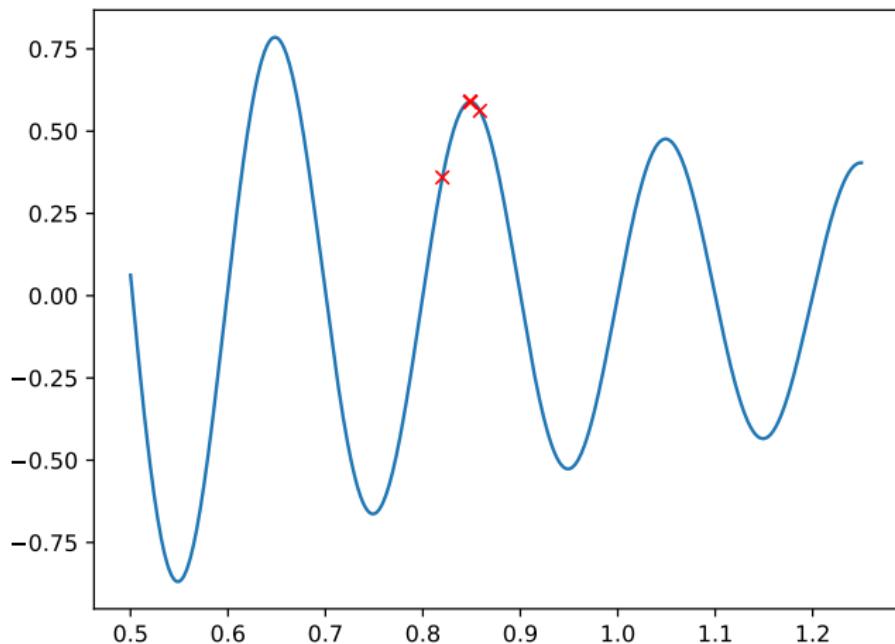
$$\nabla_{\mathbf{x}} f(\mathbf{x}_k + \Delta \mathbf{x}) \approx \nabla_{\mathbf{x}} f(\mathbf{x}_k) + \frac{\partial}{\partial \mathbf{x}} (\nabla_{\mathbf{x}} f(\mathbf{x}_k)) \Delta \mathbf{x} = 0$$

$$\nabla_{\mathbf{x}} f(\mathbf{x}_k + \Delta \mathbf{x}) \approx \nabla_{\mathbf{x}} f(\mathbf{x}_k) + \nabla_{\mathbf{xx}}^2 f(\mathbf{x}_k) \Delta \mathbf{x} = 0$$

$$\Delta \mathbf{x} = - \left(\nabla_{\mathbf{xx}}^2 f(\mathbf{x}_k) \right)^{-1} \nabla_{\mathbf{x}} f(\mathbf{x}_k)$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta \mathbf{x}$$

Optimization via Root Finding - Code Example



$$\min_{\mathbf{x}} f(\mathbf{x}), f(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}$$

Necessary Condition for solution (when f is smooth):

$$\frac{\partial f}{\partial \mathbf{x}} \Big|_{\mathbf{x}^*} = 0$$

Not sufficient! Let's have a look at the 1D case:

Optimization via Root Finding (2)

$$\min_{\mathbf{x}} f(\mathbf{x}), f(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}$$

Necessary Condition for solution (when f is smooth):

$$\frac{\partial f}{\partial \mathbf{x}} \Big|_{\mathbf{x}^*} = 0$$

Not sufficient! Let's have a look at the 1D case:

$$\Delta \mathbf{x} = \underbrace{-}_{\text{"descent"} } \underbrace{\left(\nabla_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}_k) \right)^{-1}}_{\text{"learning rate"} } \underbrace{\nabla_{\mathbf{x}} f(\mathbf{x}_k)}_{\text{"gradient"}}$$

Optimization via Root Finding (3)

Let's have a look at the 1D case:

$$\Delta \mathbf{x} = \underbrace{-}_{\text{"descent"} } \underbrace{\left(\nabla_{\mathbf{xx}}^2 f(\mathbf{x}_k) \right)^{-1}}_{\text{"learning rate"} } \underbrace{\nabla_{\mathbf{x}} f(\mathbf{x}_k)}_{\text{"gradient"}}$$

Optimization via Root Finding (3)

Let's have a look at the 1D case:

$$\Delta \mathbf{x} = \underbrace{-}_{\text{"descent"} } \underbrace{\left(\nabla_{\mathbf{xx}}^2 f(\mathbf{x}_k) \right)^{-1}}_{\text{"learning rate"} } \underbrace{\nabla_{\mathbf{x}} f(\mathbf{x}_k)}_{\text{"gradient"}}$$

- If $\nabla_{\mathbf{xx}}^2 f(\mathbf{x}) > 0$, then “**descent**”,
- If $\nabla_{\mathbf{xx}}^2 f(\mathbf{x}) < 0$, then “**ascent**”!

Optimization via Root Finding (3)

Let's have a look at the 1D case:

$$\Delta \mathbf{x} = \underbrace{-}_{\text{"descent"} } \underbrace{\left(\nabla_{\mathbf{xx}}^2 f(\mathbf{x}_k) \right)^{-1}}_{\text{"learning rate"} } \underbrace{\nabla_{\mathbf{x}} f(\mathbf{x}_k)}_{\text{"gradient"}}$$

- If $\nabla_{\mathbf{xx}}^2 f(\mathbf{x}) > 0$, then “**descent**”,
- If $\nabla_{\mathbf{xx}}^2 f(\mathbf{x}) < 0$, then “**ascent**”!

In \mathbb{R}^N :

- If $\nabla_{\mathbf{xx}}^2 f(\mathbf{x}) > 0$ (positive definite), “**descent**”,
- If $\nabla_{\mathbf{xx}}^2 f(\mathbf{x}) < 0$ (negative definite), “**ascent**”!

“Damped” Newton’s Method

So, how do we solve this?

“Damped” Newton’s Method

So, how do we solve this? Many was! Simplest yet effective is to add regularization or “damping”:

$$\mathbf{H} = \nabla_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}_k)$$

while $\mathbf{H} \leq 0$:

$$\mathbf{H} = \mathbf{H} + \beta \mathbf{I}$$

$$\Delta \mathbf{x} = -\mathbf{H}^{-1} \nabla_{\mathbf{x}} f(\mathbf{x}_k)$$

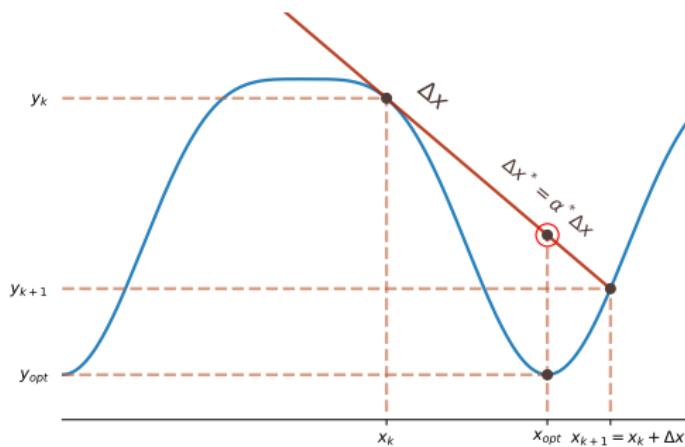
$$\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta \mathbf{x}$$

Backtracking Line Search

- Often Δx overshoots the real optimum. Why?

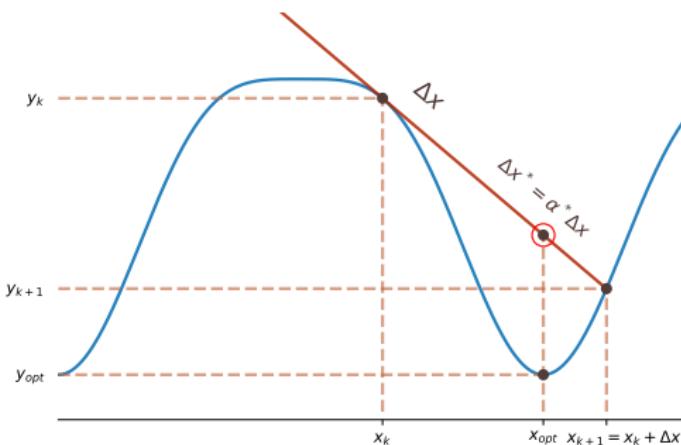
Backtracking Line Search

- Often Δx overshoots the real optimum. Why?
 - We are linearizing and taking a step on this line!!



Backtracking Line Search

- Often Δx overshoots the real optimum. Why?
 - We are linearizing and taking a step on this line!!



- To fix this issue:
 - We check if $f(\mathbf{x}_k + \Delta \mathbf{x})$ is good
 - If not, we “backtrack”; aka, we decrease $\Delta \mathbf{x}$
 - Repeat until we have a nice reduction

Armijo Rule:

$$\alpha = 1$$

while $f(\mathbf{x}_k + \alpha \Delta \mathbf{x}) > f(\mathbf{x}_k) + b\alpha \nabla_{\mathbf{x}} f(\mathbf{x}_k)^T \Delta \mathbf{x}$:

$$\alpha = c\alpha$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \Delta \mathbf{x}$$

with $0 < c < 1$, α is like a “step-size”, and b is the “tolerance”.
Usually, $c = 0.5$ and b in range from $1e^{-4}$ to 0.1.

Armijo Rule:

$$\alpha = 1$$

while $f(\mathbf{x}_k + \alpha \Delta \mathbf{x}) > f(\mathbf{x}_k) + b\alpha \nabla_{\mathbf{x}} f(\mathbf{x}_k)^T \Delta \mathbf{x}$:

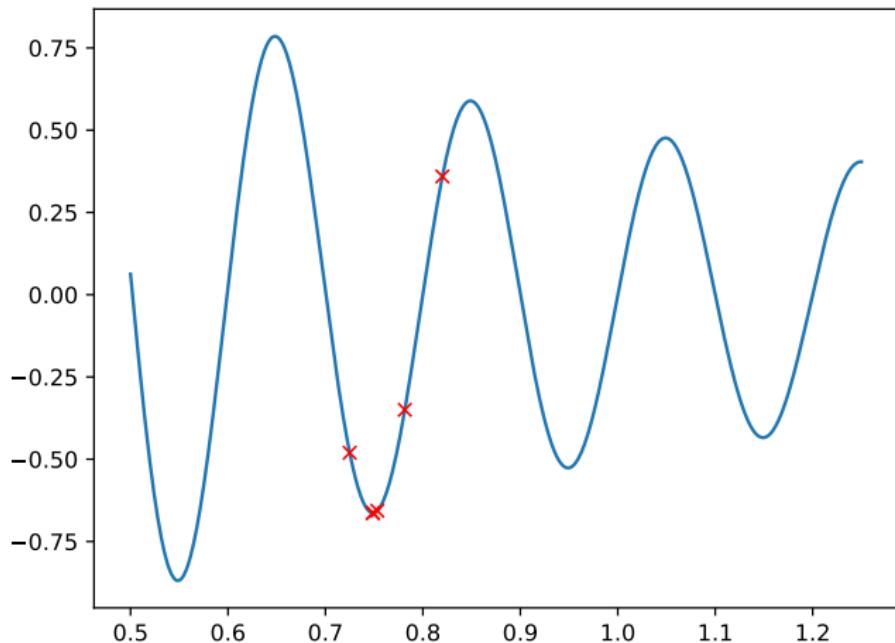
$$\alpha = c\alpha$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \Delta \mathbf{x}$$

with $0 < c < 1$, α is like a “step-size”, and b is the “tolerance”.
Usually, $c = 0.5$ and b in range from $1e^{-4}$ to 0.1.

Intuition: The actual step needs to agree with the Taylor expansion (linearization) up to a tolerance.

Newton's Method with tricks - Code Example



Newton's Method:

- Quadratic Convergence
- Very accurate
- $O(N^3)$ for solving the linear system
 - We can improve this by exploiting the structure of the problem
- **Newton's Method for Optimization:**
 - Local! It will only find the “closest” minimum point
 - **We need to “fix” it in order to optimize!**

Gradient Descent

What if we do not have the Hessian?

- We can approximate it: **Quasi-Newton Methods** (e.g. BFGS)
- We can use the gradient only!

¹Toussaint, Marc. 2012. Some Notes on Gradient Descent.

Gradient Descent

What if we do not have the Hessian?

- We can approximate it: **Quasi-Newton Methods** (e.g. BFGS)
- We can use the gradient only!

Let's go back to our update equation:

$$\Delta \mathbf{x} = \underbrace{-}_{\text{"descent"}}, \underbrace{\alpha}_{\text{"learning rate"}}, \underbrace{\nabla_{\mathbf{x}} f(\mathbf{x}_k)}_{\text{"gradient"}}$$

How to choose the “learning rate”?

¹Toussaint, Marc. 2012. Some Notes on Gradient Descent.

Gradient Descent

What if we do not have the Hessian?

- We can approximate it: **Quasi-Newton Methods** (e.g. BFGS)
- We can use the gradient only!

Let's go back to our update equation:

$$\Delta \mathbf{x} = \underbrace{-}_{\text{"descent"}}, \underbrace{\alpha}_{\text{"learning rate"}}, \underbrace{\nabla_{\mathbf{x}} f(\mathbf{x}_k)}_{\text{"gradient"}}$$

How to choose the “learning rate”?

- 1 Initialize $\alpha = \alpha_0$
- 2 $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \nabla_{\mathbf{x}} f(\mathbf{x}_k)$
- 3 if $f(\mathbf{x}_{k+1}) > f(\mathbf{x}_k)$ then $\alpha = \alpha \beta$ ($0 < \beta < 1$) and $\mathbf{x}_{k+1} = \mathbf{x}_k$
- 4 if $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$ then $\alpha = \alpha \zeta$ ($\zeta > 1$)
- 5 $k = k + 1$ and back to step 2 until we converge

It is guaranteed to always converge to a local minimum¹! Usually $\beta = 0.5, \zeta = 1.2$.

¹Toussaint, Marc. 2012. Some Notes on Gradient Descent.

Gradient Descent With Momentum

$$\Delta \mathbf{x} = \underbrace{-}_{\text{"descent"}}, \underbrace{\alpha}_{\text{"learning rate"}}, \underbrace{\nabla_{\mathbf{x}} f(\mathbf{x}_k)}_{\text{"gradient"}}$$

How to choose the “learning rate”?

- ...

Gradient Descent With Momentum

$$\Delta \mathbf{x} = \underbrace{-}_{\text{"descent"}}, \underbrace{\alpha}_{\text{"learning rate"}}, \underbrace{\nabla_{\mathbf{x}} f(\mathbf{x}_k)}_{\text{"gradient"}}$$

How to choose the “learning rate”?

- ...
- Momentum:
 - [1] Initialize $\mathbf{m}_0 = \mathbf{0}, \eta \in [0, 1]$
 - [2] $\mathbf{m}_{k+1} = \eta \mathbf{m}_k + (1 - \eta) \alpha \nabla_{\mathbf{x}} f(\mathbf{x}_k)$
 - [3] $\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{m}_{k+1}$
 - [4] $k = k + 1$ and back to step 2 until we converge

Faster convergence and “smooths out” the noise in gradient steps.

Adam Optimizer

$$\Delta \mathbf{x} = \underbrace{-}_{\text{"descent"}}, \underbrace{\alpha}_{\text{"learning rate"}}, \underbrace{\nabla_{\mathbf{x}} f(\mathbf{x}_k)}_{\text{"gradient"}}$$

How to choose the “learning rate”?

Adam Optimizer

$$\Delta \mathbf{x} = \underbrace{-}_{\text{"descent"}}, \underbrace{\alpha}_{\text{"learning rate"}}, \underbrace{\nabla_{\mathbf{x}} f(\mathbf{x}_k)}_{\text{"gradient"}}$$

How to choose the “learning rate”?

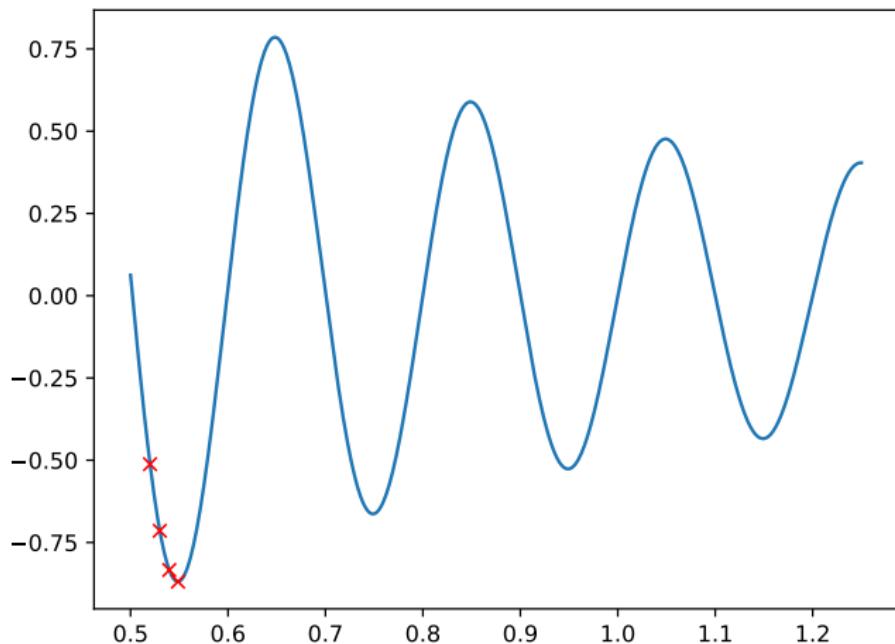
- **Adam Optimizer:**

- 1 Initialize $\mathbf{m}_0 = \mathbf{v}_0 = \mathbf{0}$, $\beta_1, \beta_2 \in [0, 1)$, $\epsilon > 0$
- 2 $\mathbf{g}_{k+1} = \nabla_{\mathbf{x}} f(\mathbf{x}_k)$
- 3 $\mathbf{m}_{k+1} = \beta_1 \mathbf{m}_k + (1 - \beta_1) \mathbf{g}_{k+1}$
- 4 $\mathbf{v}_{k+1} = \beta_2 \mathbf{v}_k + (1 - \beta_2) \mathbf{g}_{k+1}^2$
- 5 $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \frac{\frac{\mathbf{m}_{k+1}}{1 - \beta_1^{k+1}}}{\sqrt{\frac{\mathbf{v}_{k+1}}{1 - \beta_2^{k+1}}} + \epsilon}$

- 6 $k = k + 1$ and back to step 2 until we converge

One of the most versatile first order optimizers. *Not always the best choice, but surely the first one!* Usually $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$.

Gradient Descent - Code Example



Stochastic Gradient Descent

- We cannot always compute the exact $\nabla_{\mathbf{x}} f(\mathbf{x}_k)$:
 - We do not have access and we approximate it via finite differences
 - Computation limitations: too big of a dataset!
 - Computation limitations: too expensive to compute the actual one, we use an approximated version!
 - ...

Stochastic Gradient Descent

- We cannot always compute the exact $\nabla_{\mathbf{x}} f(\mathbf{x}_k)$:
 - We do not have access and we approximate it via finite differences
 - Computation limitations: too big of a dataset!
 - Computation limitations: too expensive to compute the actual one, we use an approximated version!
 - ...
- **Stochastic Gradient Descent** (SGD) is just acknowledging this fact!
- Under assumptions and mitigation measures, we can still theoretically converge to a local optimum

- **Pros:**

- Can handle large datasets
- Can handle noisy gradients
- Can handle non-smooth functions

- **Pros:**

- Can handle large datasets
- Can handle noisy gradients
- Can handle non-smooth functions

- **Cons:**

- Convergence is not guaranteed
- May oscillate around the optimum
- Requires careful tuning of hyperparameters

- **What if we cannot compute the gradient at all?**

- What if we cannot compute the gradient at all?
- We know that $\nabla_x f(\mathbf{x}) = \lim_{h \rightarrow 0} \frac{f(\mathbf{x}+h) - f(\mathbf{x})}{h}$

- What if we cannot compute the gradient at all?
- We know that $\nabla_{\mathbf{x}} f(\mathbf{x}) = \lim_{h \rightarrow 0} \frac{f(\mathbf{x}+h) - f(\mathbf{x})}{h}$
- Idea: Let's use $\frac{f(\mathbf{x}+h) - f(\mathbf{x})}{h}$ as an approximation of the gradient!
- As we make h smaller, we should be getting better approximations! Up to a limit! We have finite precision!
- In N-D, this looks like this:

$$\nabla_{\mathbf{x}} f(\mathbf{x}) \approx \begin{bmatrix} \frac{f(\mathbf{x} + h\mathbf{e}_1) - f(\mathbf{x})}{h} \\ \vdots \\ \frac{f(\mathbf{x} + h\mathbf{e}_D) - f(\mathbf{x})}{h} \end{bmatrix}$$

where \mathbf{e}_i is a direction vector (aka all zeros instead of the i -th component that is equal to 1).

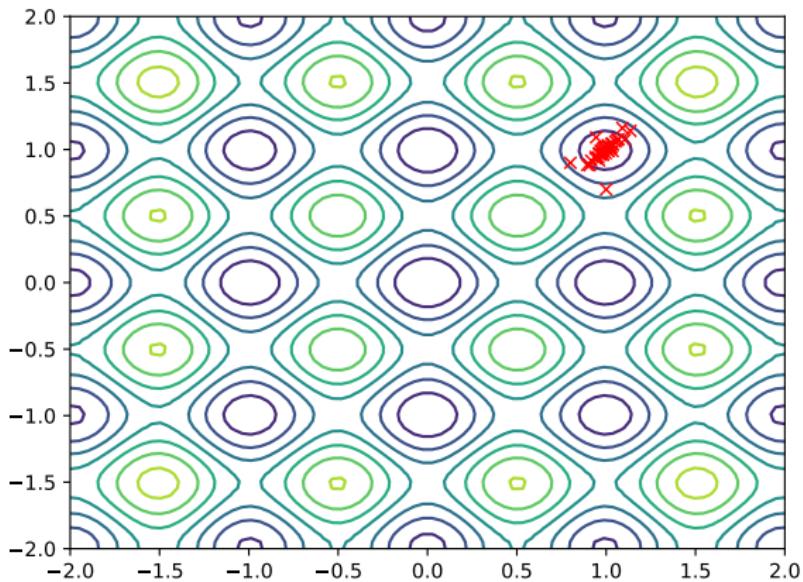
Gradient Approximation via Finite Differences (2)

- In practice we use $\frac{f(\mathbf{x}+h) - f(\mathbf{x}-h)}{2h}$
- This has an error $O(h^2)$
- Numerically more stable (aka we can use smaller h)
- In N-D, this looks like this:

$$\nabla_{\mathbf{x}} f(\mathbf{x}) \approx \begin{bmatrix} \frac{f(\mathbf{x} + h\mathbf{e}_1) - f(\mathbf{x} - h\mathbf{e}_1)}{2h} \\ \vdots \\ \frac{f(\mathbf{x} + h\mathbf{e}_D) - f(\mathbf{x} - h\mathbf{e}_D)}{2h} \end{bmatrix}$$

where \mathbf{e}_i is a direction vector (aka all zeros instead of the i -th component that is equal to 1).

Gradient Approximation - Code Example (2)



Pros:

- Simple
- Easy to implement
- Accurate results for small h

Cons:

- We need $2D$ function calls for one gradient approximation
- If your function is “expensive”, this is a bad way of approximating the gradient
- Many numerical errors can arise

- How can we use the finite difference idea and random sampling to approximate more efficiently high-dimensional gradients?
- **Simple idea:**

$$\nabla_{\mathbf{x}} f(\mathbf{x}) \approx \sum_{i=1}^N \frac{(f(\mathbf{x} + \sigma \boldsymbol{\epsilon}_i) - f(\mathbf{x} - \sigma \boldsymbol{\epsilon}_i))}{2N\sigma} \boldsymbol{\epsilon}_i$$

where $\boldsymbol{\epsilon}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

- How can we use the finite difference idea and random sampling to approximate more efficiently high-dimensional gradients?
- **Simple idea:**

$$\nabla_{\mathbf{x}} f(\mathbf{x}) \approx \sum_{i=1}^N \frac{(f(\mathbf{x} + \sigma \epsilon_i) - f(\mathbf{x} - \sigma \epsilon_i))}{2N\sigma} \epsilon_i$$

where $\epsilon_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

- **Why is this efficient?**

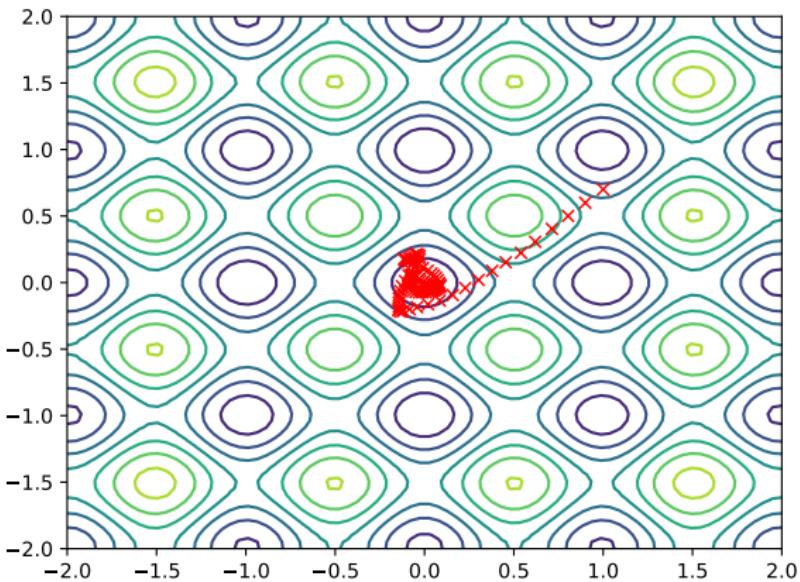
- How can we use the finite difference idea and random sampling to approximate more efficiently high-dimensional gradients?
- **Simple idea:**

$$\nabla_{\mathbf{x}} f(\mathbf{x}) \approx \sum_{i=1}^N \frac{(f(\mathbf{x} + \sigma \epsilon_i) - f(\mathbf{x} - \sigma \epsilon_i))}{2N\sigma} \epsilon_i$$

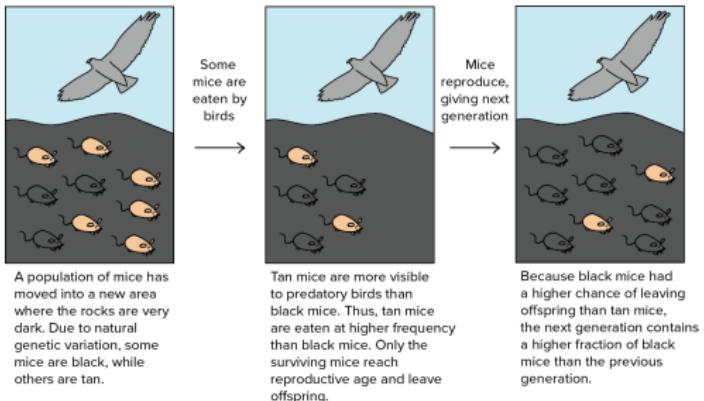
where $\epsilon_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

- Why is this efficient?
- Why is this effective?

Gradient with Random Samples - Code Example



Evolutionary Algorithms



Source: Khan Academy "Darwin, evolution, & natural selection"

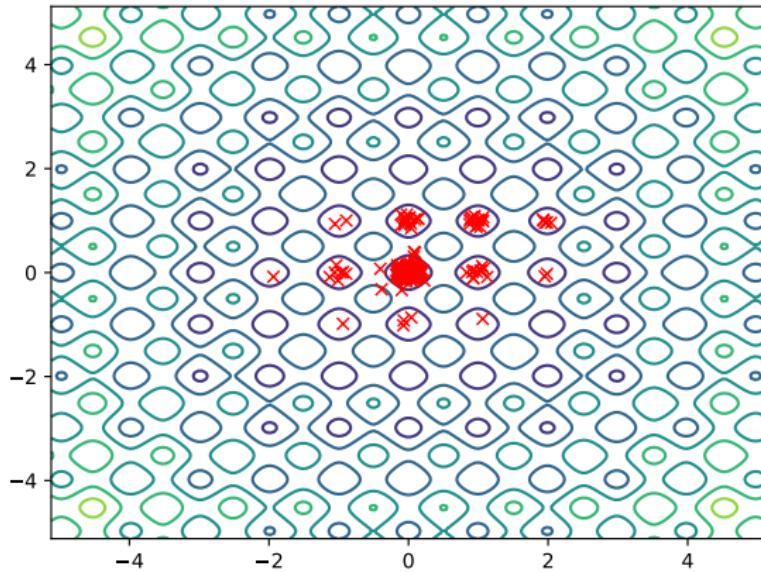
- The previous approximation is part of an *evolutionary algorithm* called **Antithetic Evolution Strategies!!**
- **Evolutionary Algorithms** (EAs) are inspired by the natural evolution of fittest.
- **Evolution Strategies** are EAs that model the population as a parameteric distribution (e.g. Gaussian) and optimize its parameters.

- 1 Generate or set the initial population
- 2 Evaluate each individual in the population
- 3 Select individuals for reproduction (*parents*)
- 4 Breed new individuals (*offspring*) or update existing individuals through crossover and mutation operations
- 5 Discard the least-fit individuals or replace them with new individuals
- 6 Possibly re-generate the population
- 7 Go back to step 2 (repeat until convergence/stopping criteria)

A Simple Genetic Algorithm

- 1 Generate or set the initial population of size N , $\mathcal{P} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- 2 Evaluate each individual in the population, $\mathcal{F} = \{f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)\}$
- 3 Sort \mathcal{F} and select the best $M < N$ individuals for reproduction,
 $\mathcal{P}_{\text{elite}} = \{\mathbf{x}_1^{\text{elite}}, \dots, \mathbf{x}_M^{\text{elite}}\}$
- 4 Set $\mathcal{O} = \emptyset$
- 5 For each $\mathbf{x}_i^{\text{elite}} \in \mathcal{P}_{\text{elite}}$:
 - 1 Select randomly another $\mathbf{x}_j^{\text{elite}} \in \mathcal{P}_{\text{elite}} (i \neq j)$
 - 2 Compute: $\mathbf{x}_i^{\text{offspring}} = \mathbf{x}_i^{\text{elite}} + \underbrace{\epsilon_1}_{\text{mutation}} + \underbrace{\epsilon_2 (\mathbf{x}_j^{\text{elite}} - \mathbf{x}_i^{\text{elite}})}_{\text{crossover}}$, where
 $\epsilon_1 \sim \mathcal{N}(\mathbf{0}, \sigma_1 I)$, $\epsilon_2 \sim \mathcal{N}(\mathbf{0}, \sigma_2 I)$.
 - 3 Add $\mathbf{x}_i^{\text{offspring}}$ to \mathcal{O}
- 6 Replace the lowest scoring individuals of \mathcal{P} with \mathcal{O}
- 7 Go back to step 2 (repeat until convergence/stopping criteria)

Genetic Algorithm - Code Example

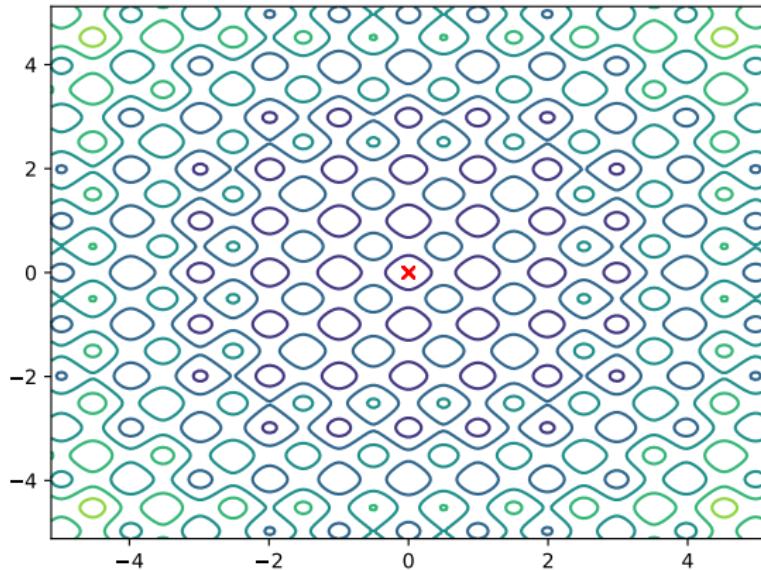


Cross Entropy Method (CEM)

Cross Entropy Method (CEM) or Simple Evolution Strategies:

- 1 The population is a Gaussian distribution: $\mathcal{N}(\mu, \text{diag}(\sigma^2))$
- 2 Generate a population of size N , $\mathcal{P} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, where $\mathbf{x}_i \sim \mathcal{N}(\mu_k, \text{diag}(\sigma_k^2))$
- 3 Evaluate each individual in the population, $\mathcal{F} = \{f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)\}$
- 4 Select the best $M < N$ individuals from \mathcal{F} : $\mathcal{P}_{\text{elite}} = \{\mathbf{x}_1^{\text{elite}}, \dots, \mathbf{x}_M^{\text{elite}}\}$
- 5 $\mu_{k+1} = \frac{1}{M} \sum_{\mathbf{x}_i^{\text{elite}} \in \mathcal{P}_{\text{elite}}} \mathbf{x}_i^{\text{elite}}$ (Sample mean)
- 6 $\sigma_{k+1}^2 = \frac{1}{M} \sum_{\mathbf{x}_i^{\text{elite}} \in \mathcal{P}_{\text{elite}}} (\mathbf{x}_i^{\text{elite}} - \mu_k)^2$ (Sample variance)
- 7 Go back to step 2 (repeat until convergence/stopping criteria)

Cross Entropy Method - Code Example



Covariance Matrix Adaptation ES (CMA-ES)

- Similar to CEM
- We use a full Covariance matrix, not just a diagonal
- Effective heuristics for adapting the Covariance matrix
- Versatile black-box optimizer!
- Very good performance with small populations!

- **Pros:**

- No need for gradients
- Can handle non-smooth functions
- Can handle noisy functions
- Global optimization

■ Pros:

- No need for gradients
- Can handle non-smooth functions
- Can handle noisy functions
- Global optimization

■ Cons:

- Slow convergence
- Requires many function evaluations
- Tricky to scale to high dimensions

$$\begin{aligned} & \min_{\mathbf{x}} f(\mathbf{x}) \\ & \text{s.t. } h(\mathbf{x}) = 0 \end{aligned}$$

where $f(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}$, $h(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}^M$.

First Order Necessary Conditions:

- 1) $\frac{\partial f}{\partial \mathbf{x}} + \boldsymbol{\lambda}^T \frac{\partial h}{\partial \mathbf{x}} = 0$
- 2) $h(\mathbf{x}) = 0$

$$\boldsymbol{\lambda} \in \mathbb{R}^M$$

$$\begin{aligned} & \min_{\mathbf{x}} f(\mathbf{x}) \\ & \text{s.t. } h(\mathbf{x}) = 0 \end{aligned}$$

where $f(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}$, $h(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}^M$.

First Order Necessary Conditions:

- 1) $\frac{\partial f}{\partial \mathbf{x}} + \boldsymbol{\lambda}^T \frac{\partial h}{\partial \mathbf{x}} = 0$
- 2) $h(\mathbf{x}) = 0$

$\boldsymbol{\lambda} \in \mathbb{R}^M$ and $\boldsymbol{\lambda}$ is called the “**Lagrange multiplier**” or “**dual variable**”.

Constrained Optimization: Equality Constraints (2)

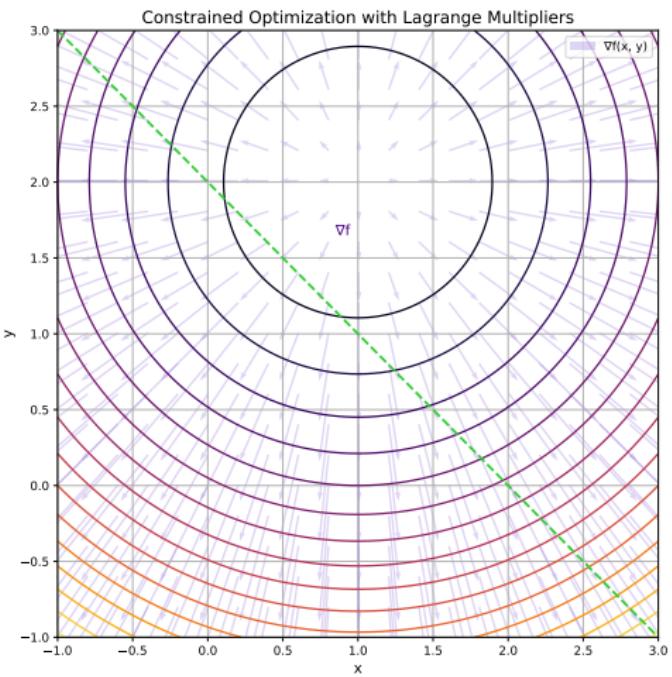
$$\min_{x,y} (x - 1)^2 + (y - 2)^2$$

$$\text{s.t. } x + y - 2 = 0$$

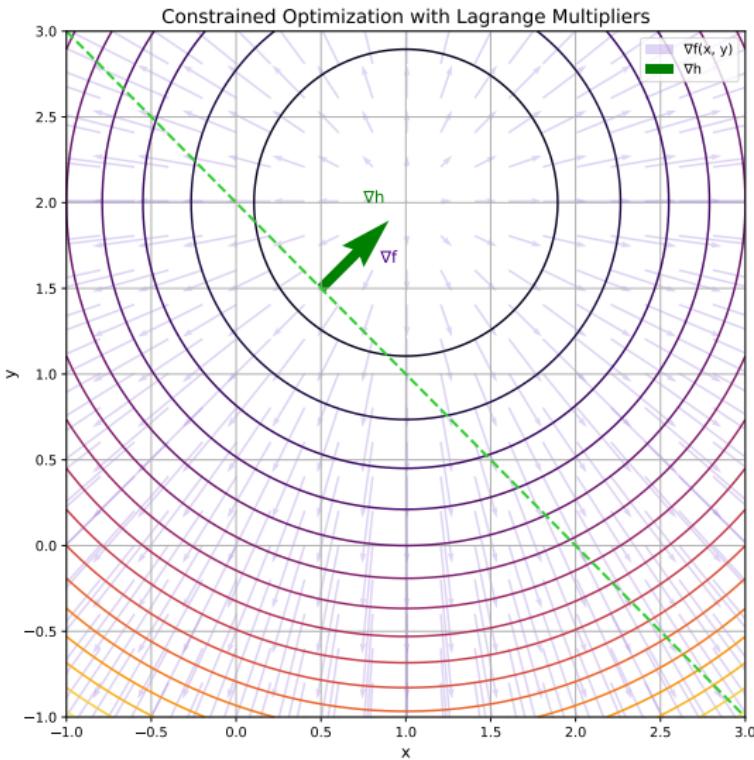
Constrained Optimization: Equality Constraints (2)

$$\min_{x,y} (x - 1)^2 + (y - 2)^2$$

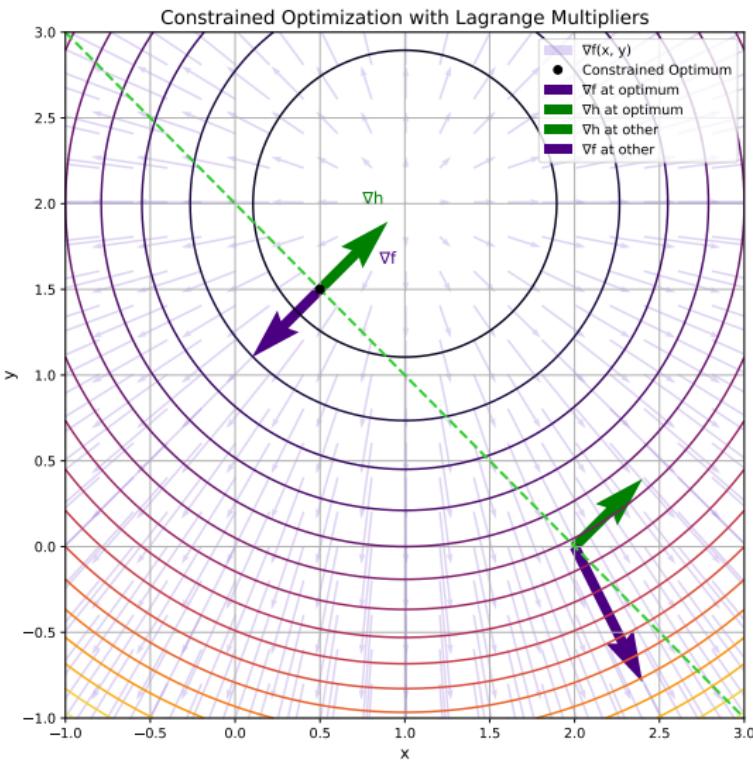
$$\text{s.t. } x + y - 2 = 0$$



Constrained Optimization: Equality Constraints (3)



Constrained Optimization: Equality Constraints (4)



Lagrangian for Equality Constraints

We define the **Lagrangian** as follows:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T h(\mathbf{x})$$

Now we can say:

- 1) $\nabla_{\mathbf{x}} \mathcal{L} = \nabla_{\mathbf{x}} f(\mathbf{x}) + \left(\frac{\partial h}{\partial \mathbf{x}} \right)^T \boldsymbol{\lambda} = 0$
- 2) $\nabla_{\boldsymbol{\lambda}} \mathcal{L} = h(\mathbf{x}) = 0$

We have recovered the original conditions! We call these conditions "**KKT Conditions**".

OK! How can we find the solution to this?

OK! How can we find the solution to this?

Idea! Let's use Newton's method!

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_k + \Delta \mathbf{x}, \boldsymbol{\lambda}_k + \Delta \boldsymbol{\lambda}) = \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_k, \boldsymbol{\lambda}_k) + \frac{\partial^2 \mathcal{L}}{\partial \mathbf{x}^2} \Big|_{\mathbf{x}_k, \boldsymbol{\lambda}_k} \Delta \mathbf{x} + \underbrace{\frac{\partial^2 \mathcal{L}}{\partial \mathbf{x} \partial \boldsymbol{\lambda}} \Big|_{\mathbf{x}_k, \boldsymbol{\lambda}_k} \Delta \boldsymbol{\lambda}}_{\left(\frac{\partial h}{\partial \mathbf{x}} \Big|_{\mathbf{x}_k} \right)^T} = 0$$

$$\nabla_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{x}_k + \Delta \mathbf{x}, \boldsymbol{\lambda}_k + \Delta \boldsymbol{\lambda}) = h(\mathbf{x}_k) + \frac{\partial h}{\partial \boldsymbol{\lambda}} \Big|_{\mathbf{x}_k} \Delta \boldsymbol{\lambda} = 0$$

Optimization with Equality Constraints (2)

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_k + \Delta \mathbf{x}, \boldsymbol{\lambda}_k + \Delta \boldsymbol{\lambda}) = \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_k, \boldsymbol{\lambda}_k) + \frac{\partial^2 \mathcal{L}}{\partial \mathbf{x}^2} \Big|_{\mathbf{x}_k, \boldsymbol{\lambda}_k} \Delta \mathbf{x} + \left(\frac{\partial h}{\partial \mathbf{x}} \Big|_{\mathbf{x}_k} \right)^T \Delta \boldsymbol{\lambda} = 0$$

$$\nabla_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{x}_k + \Delta \mathbf{x}, \boldsymbol{\lambda}_k + \Delta \boldsymbol{\lambda}) = h(\mathbf{x}_k) + \frac{\partial h}{\partial \mathbf{x}} \Big|_{\mathbf{x}_k} \Delta \mathbf{x} = 0$$

This is a linear system in $\mathbf{z} = \begin{bmatrix} \Delta \mathbf{x} \\ \Delta \boldsymbol{\lambda} \end{bmatrix}$:

$$\begin{bmatrix} \frac{\partial^2 \mathcal{L}}{\partial \mathbf{x}^2} \Big|_{\mathbf{x}_k, \boldsymbol{\lambda}_k} & \left(\frac{\partial h}{\partial \mathbf{x}} \Big|_{\mathbf{x}_k} \right)^T \\ \frac{\partial h}{\partial \mathbf{x}} \Big|_{\mathbf{x}_k} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \Delta \mathbf{x} \\ \Delta \boldsymbol{\lambda} \end{bmatrix} = \begin{bmatrix} -\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_k, \boldsymbol{\lambda}_k) \\ -h(\mathbf{x}_k) \end{bmatrix}$$

This is often called the “**KKT System**”.

$$\frac{\partial^2 \mathcal{L}}{\partial \mathbf{x}^2} = \nabla_{\mathbf{x}\mathbf{x}}^2 f + \frac{\partial}{\partial \mathbf{x}} \left[\left(\frac{\partial h}{\partial \mathbf{x}} \right)^T \boldsymbol{\lambda} \right]$$

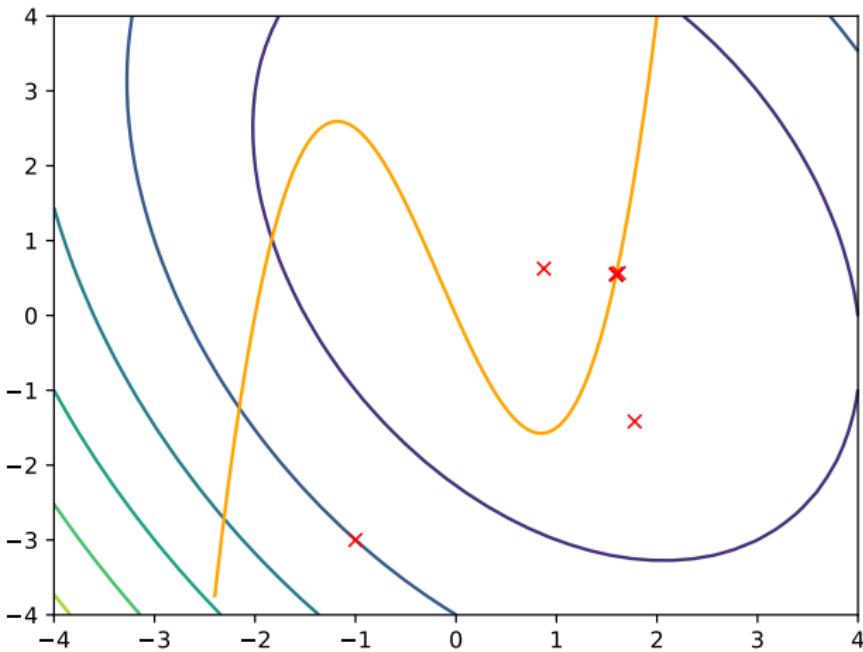
The second term here is:

- A third rank tensor
- Expensive to compute (there are “complicated” ways to compute this faster, e.g. pullback and auto-diff)
- Can be “ugly”/ill-conditioned

In practice, we often remove this term:

- The algorithm is now called “Gauss-Newton”
- In theory it has slower convergence
- In practice it takes more iterations, but each iteration is faster!

Gauss-Netwon Method - Code Example



$$\begin{aligned} & \min_{\mathbf{x}} f(\mathbf{x}) \\ & \text{s.t. } g(\mathbf{x}) \leq 0 \end{aligned}$$

where $f(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}$, $g(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}^M$. We define the **Lagrangian** as follows:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\mu}) = f(\mathbf{x}) + \boldsymbol{\mu}^T g(\mathbf{x})$$

First Order Necessary (KKT) Conditions:

- 1) $\nabla_{\mathbf{x}} f(\mathbf{x}) + \left(\frac{\partial g}{\partial \mathbf{x}} \right)^T \boldsymbol{\mu} = 0$, “stationarity”
- 2) $g(\mathbf{x}) \leq 0$, “primal feasibility”
- 3) $\boldsymbol{\mu} \geq 0$, “dual feasibility”
- 4) $\boldsymbol{\mu}^T g(\mathbf{x}) = 0$, “complementarity”

We cannot solve the above in closed form as in the equality case!!

Many methods to solve this:

- Augmented Lagrangian Method (ALM)
- Sequential Quadratic Programming (SQP)
- Interior Point Methods
- ...

We define the **Augmented Lagrangian** as:

$$\mathcal{L}_\rho(\mathbf{x}, \boldsymbol{\mu}) = f(\mathbf{x}) + \boldsymbol{\mu}^T g(\mathbf{x}) + \frac{\rho}{2} \|\max(0, g(\mathbf{x}))\|^2$$

$\frac{\rho}{2} \|\max(0, g(\mathbf{x}))\|^2$ is called the **penalty** term.

Pseudocode:

$$1) \quad \mathbf{x}_{k+1} = \min_{\mathbf{x}} \mathcal{L}_\rho(\mathbf{x}, \boldsymbol{\mu}_k)$$

$$2) \quad \boldsymbol{\mu}_{k+1} = \max\left(\mathbf{0}, \boldsymbol{\mu}_k + \rho g(\mathbf{x}_{k+1})\right)$$

$$3) \quad \rho = \alpha \rho, \text{ with } \alpha \in \mathbb{R}^+$$

Usually $\alpha \approx 10$ and we use (2) only when the constraints are violated. The method comes with strong convergence guarantees.

How can we solve (1)?

Augmented Lagrangian Method - Minimizing \mathcal{L}_ρ

- We need to solve $\min_{\mathbf{x}} \mathcal{L}_\rho(\mathbf{x}, \boldsymbol{\mu})$
- Important to note that $\boldsymbol{\mu}$ and ρ are constants!

Augmented Lagrangian Method - Minimizing \mathcal{L}_ρ

- We need to solve $\min_{\mathbf{x}} \mathcal{L}_\rho(\mathbf{x}, \boldsymbol{\mu})$
- Important to note that $\boldsymbol{\mu}$ and ρ are constants!
- Let's apply Newton's Method:

$$\nabla_{\mathbf{x}} \mathcal{L}_\rho(\mathbf{x}_k + \Delta \mathbf{x}, \boldsymbol{\mu}_k) \approx$$

$$\nabla_{\mathbf{x}} \mathcal{L}_\rho(\mathbf{x}_k, \boldsymbol{\mu}_k) + \frac{\partial}{\partial \mathbf{x}} \left(\nabla_{\mathbf{x}} \mathcal{L}_\rho(\mathbf{x}_k, \boldsymbol{\mu}_k) \right) \Delta \mathbf{x} = 0$$

- We need the gradient of \mathcal{L}_ρ :

$$\begin{aligned}\nabla_{\mathbf{x}} \mathcal{L}_\rho((\mathbf{x}, \boldsymbol{\mu})) &= \left(\frac{\partial f}{\partial \mathbf{x}} + \boldsymbol{\mu}^T \frac{\partial g}{\partial \mathbf{x}} + \rho g(\mathbf{x})^T \frac{\partial g}{\partial \mathbf{x}} \right)^T \\ &= \left(\frac{\partial f}{\partial \mathbf{x}} + \left(\boldsymbol{\mu} + \rho g(\mathbf{x}) \right)^T \frac{\partial g}{\partial \mathbf{x}} \right)^T = 0\end{aligned}$$

Augmented Lagrangian Method - Minimizing \mathcal{L}_ρ (2)

$$\nabla_{\mathbf{x}} \mathcal{L}_\rho(\mathbf{x}_k + \Delta \mathbf{x}, \boldsymbol{\mu}_k) \approx$$

$$\nabla_{\mathbf{x}} \mathcal{L}_\rho(\mathbf{x}_k, \boldsymbol{\mu}_k) + \frac{\partial}{\partial \mathbf{x}} \left(\nabla_{\mathbf{x}} \mathcal{L}_\rho(\mathbf{x}_k, \boldsymbol{\mu}_k) \right) \Delta \mathbf{x} = 0$$

We can use the Gauss-Newton version and:

$$\frac{\partial}{\partial \mathbf{x}} \left(\nabla_{\mathbf{x}} \mathcal{L}_\rho(\mathbf{x}_k, \boldsymbol{\mu}_k) \right) \approx \nabla_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}_k) + \rho \left(\frac{\partial g}{\partial \mathbf{x}} \Big|_{\mathbf{x}_k} \right)^T \frac{\partial g}{\partial \mathbf{x}} \Big|_{\mathbf{x}_k}$$

$$\Delta \mathbf{x} = - \left(\nabla_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}_k) + \rho \left(\frac{\partial g}{\partial \mathbf{x}} \Big|_{\mathbf{x}_k} \right)^T \frac{\partial g}{\partial \mathbf{x}} \Big|_{\mathbf{x}_k} \right)^{-1} \nabla_{\mathbf{x}} \mathcal{L}_\rho(\mathbf{x}_k, \boldsymbol{\mu}_k)$$

Merit Functions

- How can we apply line search in the constrained case?

- How can we apply line search in the constrained case?
- We define a **merit function**, $P(\mathbf{x})$
- Then we apply the Armijo rule to it

$$\alpha = 1$$

while $P(\mathbf{x}_k + \alpha \Delta \mathbf{x}) > P(\mathbf{x}_k) + b\alpha \nabla_{\mathbf{x}} P(\mathbf{x}_k)^T \Delta \mathbf{x}$:

$$\alpha = c\alpha$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \Delta \mathbf{x}$$

Full constrained minimization:

$$\begin{aligned} & \min_{\mathbf{x}} f(\mathbf{x}) \\ \text{s.t. } & g(\mathbf{x}) \leq 0 \\ & h(\mathbf{x}) = 0 \end{aligned}$$

The **Lagrangian** is given:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T h(\mathbf{x}) + \boldsymbol{\mu}^T g(\mathbf{x})$$

First Order Necessary (KKT) Conditions:

- 1) $\nabla_{\mathbf{x}} f(\mathbf{x}) + \left(\frac{\partial h}{\partial \mathbf{x}} \right)^T \boldsymbol{\lambda} + \left(\frac{\partial g}{\partial \mathbf{x}} \right)^T \boldsymbol{\mu} = 0,$ “stationarity”
- 2) $g(\mathbf{x}) \leq 0$ and $h(\mathbf{x}) = 0,$ “primal feasibility”
- 3) $\boldsymbol{\mu} \geq 0,$ “dual feasibility”
- 4) $\boldsymbol{\mu}^T g(\mathbf{x}) = 0,$ “complementarity”

Usual Merit Functions

- KKT residual

$$P(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \frac{1}{2} \left\| \begin{bmatrix} \nabla_{\mathbf{x}} \mathcal{L} \\ h(\mathbf{x}) \\ \min(0, g(\mathbf{x})) \end{bmatrix} \right\|^2$$

- Weighted Sum

$$P(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \frac{1}{2} w \left\| \begin{bmatrix} h(\mathbf{x}) \\ \min(0, g(\mathbf{x})) \end{bmatrix} \right\|^2$$

- Augmented Lagrangian

$$P(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathcal{L}_{\rho}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$$

Quadratic Programming (QP)

$$\begin{aligned}\min_{\mathbf{x}} f(\mathbf{x}) &= \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{q}^T \mathbf{x} \\ \text{s.t. } \mathbf{A} \mathbf{x} - \mathbf{b} &= \mathbf{0} \\ \mathbf{C} \mathbf{x} - \mathbf{d} &\leq \mathbf{0}\end{aligned}$$

where $\mathbf{x}, \mathbf{q} \in \mathbb{R}^N$, $\mathbf{Q} > 0 \in \mathbb{R}^{N \times N}$. Let's solve this with the **Augmented Lagrangian Method**:

Quadratic Programming (QP)

$$\begin{aligned}\min_{\mathbf{x}} f(\mathbf{x}) &= \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{q}^T \mathbf{x} \\ \text{s.t. } \mathbf{A} \mathbf{x} - \mathbf{b} &= \mathbf{0} \\ \mathbf{C} \mathbf{x} - \mathbf{d} &\leq \mathbf{0}\end{aligned}$$

where $\mathbf{x}, \mathbf{q} \in \mathbb{R}^N$, $\mathbf{Q} > 0 \in \mathbb{R}^{N \times N}$. Let's solve this with the **Augmented Lagrangian Method**:

$$\begin{aligned}\mathcal{L}_\rho(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) &= f(\mathbf{x}) + \boldsymbol{\lambda}^T (\mathbf{A} \mathbf{x} - \mathbf{b}) + \boldsymbol{\mu}^T (\mathbf{C} \mathbf{x} - \mathbf{d}) \\ &\quad + \frac{\rho}{2} \|\mathbf{A} \mathbf{x} - \mathbf{b}\|^2 + \frac{\rho}{2} \|\max(\mathbf{0}, \mathbf{C} \mathbf{x} - \mathbf{d})\|^2\end{aligned}$$

KKT Conditions:

- 1) $\mathbf{Q} \mathbf{x} + \mathbf{q} + \mathbf{A}^T \boldsymbol{\lambda} + \mathbf{C}^T \boldsymbol{\mu} = \mathbf{0}$
- 2) $\mathbf{A} \mathbf{x} - \mathbf{b} = \mathbf{0}$ and $\mathbf{C} \mathbf{x} - \mathbf{d} \leq \mathbf{0}$
- 3) $\boldsymbol{\mu} \geq \mathbf{0}$
- 4) $\boldsymbol{\mu}^T (\mathbf{C} \mathbf{x} - \mathbf{d}) = \mathbf{0}$

What other methods are there?

- **Active Set Methods:**

- Try to guess which constraints are active (equal to zero) or not
- Solve equality constrained optimization only with active constraints
- Need good heuristic for identifying potentially active constraints

- **Barrier/Interior Point (IP) Methods:**

- Add inequality constraints as “barriers” /penalties in the cost function. For example, $f'(\mathbf{x}) = f(\mathbf{x}) - \frac{1}{\rho} \log(-g(\mathbf{x}))$
- Very good for convex problems

What other methods are there?

- **Penalty Methods:**

- Add constraint violation into the cost function. For example,
$$f'(\mathbf{x}) = f(\mathbf{x}) + \frac{\rho}{2} \|\max(\mathbf{0}, g(\mathbf{x}))\|^2$$
- Similar to IP methods, but there is a key difference: IP methods blow up at the boundary, penalty methods start pushing only if there is a violation!
- Can easily get ill-conditioned
- Cannot get good accuracy

- **Sequential Quadratic Programming (SQP):**

- Sequential QP problems
- General idea:
 - Fit a quadratic approximation of the cost function
 - Linearize the constraints
 - Solve the produced QP with slack variables for feasibility
 - Iterate until convergence

Final Verdict on Optimization

- Newton's Method is very effective; Gauss-Newton the one used mostly in practice
- Transforming your problem into well known forms (e.g. QP) can greatly help
- Stochastic Gradient Descent is very effective for large scale or noisy problems
- Evolutionary Algorithms are very effective for non-smooth problems, but do not scale easily on large scale problems
- A lot of good libraries

Continuous Time:

$$\begin{aligned} \underset{\mathbf{x}(t), \mathbf{u}(t)}{\operatorname{argmin}} \mathcal{J}(\mathbf{x}(t), \mathbf{u}(t)) &= \int_{t_0}^{t_f} \ell(\mathbf{x}(t), \mathbf{u}(t)) dt + \ell_F(\mathbf{x}(t_f)) \\ \text{s.t. } \dot{\mathbf{x}}(t) &= f(\mathbf{x}(t), \mathbf{u}(t)) \end{aligned}$$

where

- $\mathbf{x}(t) \in \mathbb{R}^N$, $\mathbf{u}(t) \in \mathbb{R}^M$ are the state and control trajectories
- $\mathcal{J}(\mathbf{x}(t), \mathbf{u}(t))$ is the “cost function”
- $\ell(\mathbf{x}(t), \mathbf{u}(t))$ is the “stage cost”
- $\ell_F(\mathbf{x}(t_f))$ is the “terminal cost”
- $\dot{\mathbf{x}}(t) = f(\mathbf{x}(t), \mathbf{u}(t))$ are the “dynamics constraints”
- We can potentially add more constraints (e.g. torque limits)

- The continuous version is “**infinite dimensional**”

- The continuous version is “**infinite dimensional**”
- The solution is open loop control trajectories!
- Only very few problems can be solved analytically

- The continuous version is “**infinite dimensional**”
- The solution is open loop control trajectories!
- Only very few problems can be solved analytically
- Let’s discretize!

Discrete Time:

$$\begin{aligned} \operatorname{argmin}_{\boldsymbol{x}_{1:K}, \boldsymbol{u}_{1:K-1}} \mathcal{J}(\boldsymbol{x}_{1:K}, \boldsymbol{u}_{1:K-1}) &= \sum_{k=1}^{K-1} \ell(\boldsymbol{x}_k, \boldsymbol{u}_k) + \ell_F(\boldsymbol{x}_K) \\ \text{s.t. } \boldsymbol{x}_{k+1} &= f_{\text{discrete}}(\boldsymbol{x}_k, \boldsymbol{u}_k) \end{aligned}$$

- $\boldsymbol{x}_k \in \mathbb{R}^N$ and $\boldsymbol{u}_k \in \mathbb{R}^M$ are vectors
- The is now “**finite dimensional**”
- The solution is open loop control trajectories!
- We usually call $\boldsymbol{x}_k, \boldsymbol{u}_k$ “**knot points**”

What is the Linear Quadratic Regulator (LQR) problem?

What is the Linear Quadratic Regulator (LQR) problem?

$$\underset{\mathbf{x}_{1:K}, \mathbf{u}_{1:K-1}}{\operatorname{argmin}} \mathcal{J}(\mathbf{x}_{1:K}, \mathbf{u}_{1:K-1}) = \sum_{k=1}^{K-1} \left(\frac{1}{2} \mathbf{x}_k^T \mathbf{Q}_k \mathbf{x}_k + \frac{1}{2} \mathbf{u}_k^T \mathbf{R}_k \mathbf{u}_k \right) + \frac{1}{2} \mathbf{x}_K^T \mathbf{Q}_K \mathbf{x}_K$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{A}_k \mathbf{x}_k + \mathbf{B}_k \mathbf{u}_k$$

$$\mathbf{Q}_k \succeq 0$$

$$\mathbf{R}_k > 0$$

- Widely used in many real applications
- The “workhorse” of optimal control
- We know “everything” about it!
- Infinite variations and extensions!
- **Time Invariant** if: $\mathbf{A}_k = \mathbf{A}$, $\mathbf{B}_k = \mathbf{B}$, $\mathbf{Q}_k = \mathbf{Q}$, $\mathbf{R}_k = \mathbf{R}$, $\forall k$

LQR as a Quadratic Programming Problem

Looking at the LQR problem, it really looks like a QP. Can we write it like one?

Looking at the LQR problem, it really looks like a QP. Can we write it like one? Yes we can! We assume that x_1 (initial conditions) is given, and define:

$$\mathbf{z} = \begin{bmatrix} \mathbf{u}_1 \\ \mathbf{x}_2 \\ \mathbf{u}_2 \\ \vdots \\ \mathbf{x}_K \end{bmatrix}, \mathbf{H} = \begin{bmatrix} \mathbf{R}_1 & & & & \\ & \mathbf{Q}_2 & & & \\ & & \mathbf{R}_2 & & \\ & & & \ddots & \\ & & & & \mathbf{Q}_K \end{bmatrix}$$

Now we can define:

$$\operatorname{argmin}_{\mathbf{z}} \mathcal{J}(\mathbf{z}) = \frac{1}{2} \mathbf{z}^T \mathbf{H} \mathbf{z}$$

s.t. “dynamics constraints”

LQR as a Quadratic Programming Problem (2)

For the dynamics constraints we have:

$$\underbrace{\begin{bmatrix} \mathbf{B}_1 & -\mathbf{I} & \mathbf{0} & \dots \\ \mathbf{0} & \mathbf{A}_2 & \mathbf{B}_2 & -\mathbf{I} & \mathbf{0} & \dots \\ & & & \ddots & & \\ & & & & \mathbf{A}_{K-1} & \mathbf{B}_{K-1} & -\mathbf{I} \end{bmatrix}}_G \begin{bmatrix} \mathbf{u}_1 \\ \mathbf{x}_2 \\ \mathbf{u}_2 \\ \vdots \\ \mathbf{x}_K \end{bmatrix} = \underbrace{\begin{bmatrix} -\mathbf{A}_1 \mathbf{x}_1 \\ \mathbf{0} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{bmatrix}}_d$$

Now we have a full QP:

$$\begin{aligned} \operatorname{argmin}_{\mathbf{z}} \mathcal{J}(\mathbf{z}) &= \frac{1}{2} \mathbf{z}^T \mathbf{H} \mathbf{z} \\ \text{s.t.} \quad \mathbf{G} \mathbf{z} - \mathbf{d} &= \mathbf{0} \end{aligned}$$

LQR as a Quadratic Programming Problem (2)

$$\begin{aligned}\operatorname{argmin}_z \mathcal{J}(z) &= \frac{1}{2} z^T H z \\ \text{s.t. } Gz - d &= 0\end{aligned}$$

Can we solve this?

LQR as a Quadratic Programming Problem (2)

$$\underset{\mathbf{z}}{\operatorname{argmin}} \mathcal{J}(\mathbf{z}) = \frac{1}{2} \mathbf{z}^T \mathbf{H} \mathbf{z}$$

s.t. $\mathbf{Gz} - \mathbf{d} = \mathbf{0}$

Can we solve this? Of course we can! The Lagrangian is:

$$\mathcal{L}(\mathbf{z}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{z}^T \mathbf{H} \mathbf{z} + \boldsymbol{\lambda}^T (\mathbf{Gz} - \mathbf{d})$$

LQR as a Quadratic Programming Problem (2)

$$\underset{\mathbf{z}}{\operatorname{argmin}} \mathcal{J}(\mathbf{z}) = \frac{1}{2} \mathbf{z}^T \mathbf{H} \mathbf{z}$$
$$\text{s.t.} \quad \mathbf{G} \mathbf{z} - \mathbf{d} = \mathbf{0}$$

Can we solve this? Of course we can! The Lagrangian is:

$$\mathcal{L}(\mathbf{z}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{z}^T \mathbf{H} \mathbf{z} + \boldsymbol{\lambda}^T (\mathbf{G} \mathbf{z} - \mathbf{d})$$

KKT Conditions:

$$\nabla_{\mathbf{z}} \mathcal{L} = \mathbf{H} \mathbf{z} + \mathbf{G}^T \boldsymbol{\lambda} = \mathbf{0}$$

$$\nabla_{\boldsymbol{\lambda}} \mathcal{L} = \mathbf{G} \mathbf{z} - \mathbf{d} = \mathbf{0}$$

LQR as a Quadratic Programming Problem (2)

$$\underset{\mathbf{z}}{\operatorname{argmin}} \mathcal{J}(\mathbf{z}) = \frac{1}{2} \mathbf{z}^T \mathbf{H} \mathbf{z}$$
$$\text{s.t.} \quad \mathbf{G} \mathbf{z} - \mathbf{d} = \mathbf{0}$$

Can we solve this? Of course we can! The Lagrangian is:

$$\mathcal{L}(\mathbf{z}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{z}^T \mathbf{H} \mathbf{z} + \boldsymbol{\lambda}^T (\mathbf{G} \mathbf{z} - \mathbf{d})$$

KKT Conditions:

$$\nabla_{\mathbf{z}} \mathcal{L} = \mathbf{H} \mathbf{z} + \mathbf{G}^T \boldsymbol{\lambda} = \mathbf{0}$$

$$\nabla_{\boldsymbol{\lambda}} \mathcal{L} = \mathbf{G} \mathbf{z} - \mathbf{d} = \mathbf{0}$$

KKT System:

$$\begin{bmatrix} \mathbf{H} & \mathbf{G}^T \\ \mathbf{G} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{z} \\ \boldsymbol{\lambda} \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{d} \end{bmatrix}$$

LQR as a Quadratic Programming Problem (2)

$$\underset{\mathbf{z}}{\operatorname{argmin}} \mathcal{J}(\mathbf{z}) = \frac{1}{2} \mathbf{z}^T \mathbf{H} \mathbf{z}$$
$$\text{s.t.} \quad \mathbf{G} \mathbf{z} - \mathbf{d} = \mathbf{0}$$

Can we solve this? Of course we can! The Lagrangian is:

$$\mathcal{L}(\mathbf{z}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{z}^T \mathbf{H} \mathbf{z} + \boldsymbol{\lambda}^T (\mathbf{G} \mathbf{z} - \mathbf{d})$$

KKT Conditions:

$$\nabla_{\mathbf{z}} \mathcal{L} = \mathbf{H} \mathbf{z} + \mathbf{G}^T \boldsymbol{\lambda} = \mathbf{0}$$

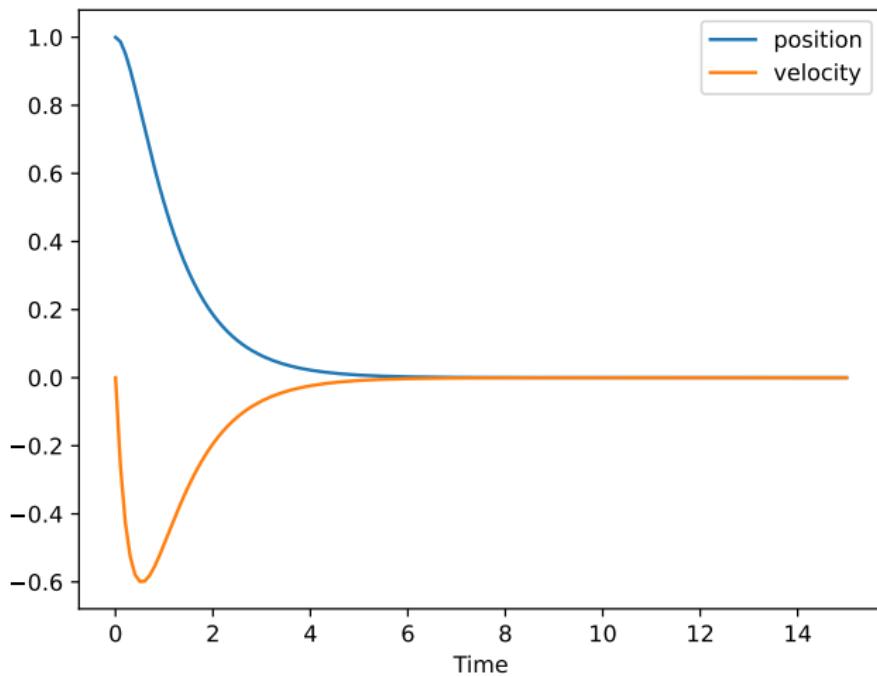
$$\nabla_{\boldsymbol{\lambda}} \mathcal{L} = \mathbf{G} \mathbf{z} - \mathbf{d} = \mathbf{0}$$

KKT System:

$$\begin{bmatrix} \mathbf{H} & \mathbf{G}^T \\ \mathbf{G} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{z} \\ \boldsymbol{\lambda} \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{d} \end{bmatrix}$$

We have the exact solution by solving ONE linear system!

LQR as a QP - Code Example



LQR via Riccati Recursion

Let's write the KKT System in "detail" (for $K = 4$):

$$\left[\begin{array}{ccccccccc} R_1 & 0 & \cdots & 0 & B_1^T & \cdots & 0 \\ 0 & Q_2 & 0 & \cdots & 0 & -I & A_2^T & \vdots \\ \vdots & 0 & R_2 & 0 & \cdots & 0 & B_2^T \\ \vdots & 0 & Q_3 & 0 & \cdots & -I & A_3^T \\ \vdots & 0 & R_3 & 0 & \cdots & B_3^T \\ 0 & 0 & \vdots & 0 & Q_4 & & -I \\ \hline B_1 & -I & 0 & 0 & \vdots & 0 & \ddots & \vdots \\ \vdots & A_2 & B_2 & -I & 0 & \vdots \\ 0 & \cdots & & A_3 & B_3 & -I & \cdots & 0 \end{array} \right] \begin{bmatrix} u_1 \\ x_2 \\ u_2 \\ x_3 \\ u_3 \\ x_4 \\ \lambda_2 \\ \lambda_3 \\ \lambda_4 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ -A_1 x_1 \\ 0 \\ 0 \end{bmatrix}$$

LQR via Riccati Recursion (2)

We first solve the last line of the upper block:

$$\mathbf{Q}_4 \mathbf{x}_4 - \lambda_4 = \mathbf{0} \Rightarrow \lambda_4 = \mathbf{Q}_K \mathbf{x}_4$$

Then we move one line up:

$$\begin{aligned}\mathbf{R}_3 \mathbf{u}_3 + \mathbf{B}_3^T \lambda_4 &= \mathbf{R}_3 \mathbf{u}_3 + \mathbf{B}_3^T \mathbf{Q}_K \mathbf{x}_4 = \mathbf{0} \\ \Rightarrow \mathbf{R}_3 \mathbf{u}_3 + \mathbf{B}_3^T \mathbf{Q}_K (\mathbf{A}_3 \mathbf{x}_3 + \mathbf{B}_3 \mathbf{u}_3) &= \mathbf{0} \\ \Rightarrow \mathbf{u}_3 &= - \underbrace{(\mathbf{R}_3 + \mathbf{B}_3^T \mathbf{Q}_K \mathbf{B}_3)^{-1} \mathbf{B}_3^T \mathbf{Q}_K \mathbf{A}_3}_{\mathbf{K}_3} \mathbf{x}_3\end{aligned}$$

One more line:

$$\begin{aligned}\mathbf{Q}_3 \mathbf{x}_3 - \lambda_3 + \mathbf{A}_3^T \lambda_4 &= \mathbf{0} \Rightarrow \mathbf{Q}_3 \mathbf{x}_3 - \lambda_3 + \mathbf{A}_3^T \mathbf{Q}_K \mathbf{x}_4 = \mathbf{0} \\ \Rightarrow \mathbf{Q}_3 \mathbf{x}_3 - \lambda_3 + \mathbf{A}_3^T \mathbf{Q}_K (\mathbf{A}_3 \mathbf{x}_3 + \mathbf{B}_3 \mathbf{u}_3) &= \mathbf{0} \\ \Rightarrow \lambda_3 &= \underbrace{\left(\mathbf{Q}_3 + \mathbf{A}_3^T \mathbf{Q}_K (\mathbf{A}_3 - \mathbf{B}_3 \mathbf{K}_3) \right)}_{\mathbf{P}_3} \mathbf{x}_3\end{aligned}$$

The recursion is now revealed:

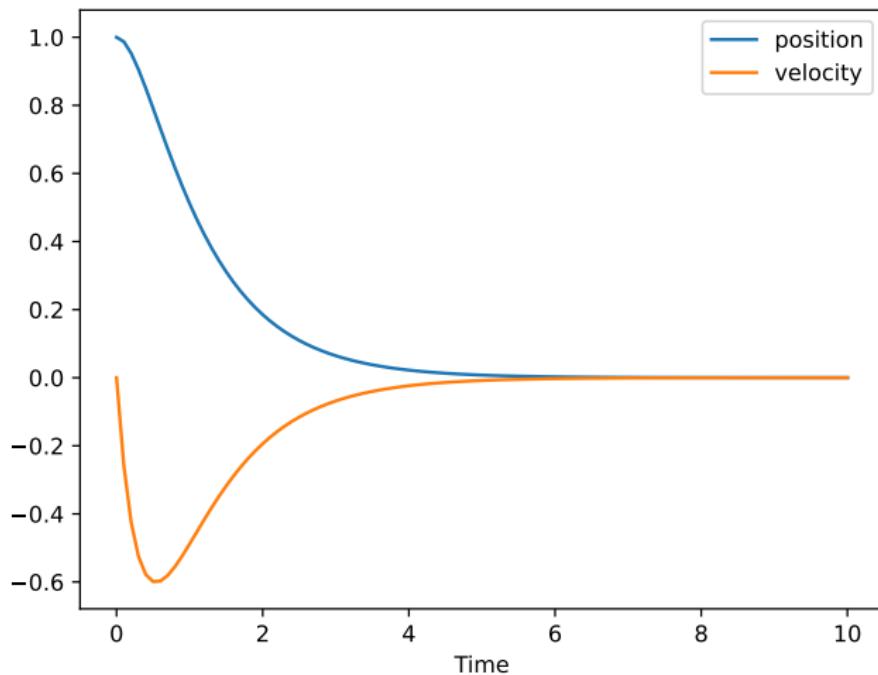
$$\mathbf{P}_K = \mathbf{Q}_K$$

$$\mathbf{K}_k = (\mathbf{R}_k + \mathbf{B}_k^T \mathbf{P}_{k+1} \mathbf{B}_k)^{-1} \mathbf{B}_k^T \mathbf{P}_{k+1} \mathbf{A}_k$$

$$\mathbf{P}_k = \mathbf{Q}_k + \mathbf{A}_k^T \mathbf{P}_{k+1} (\mathbf{A}_k - \mathbf{B}_k \mathbf{K}_k)$$

- QP has complexity $O(K^3(N + M)^3)$
- Riccati recursion has complexity $O(K(N + M)^3)$
- The Riccati recursion is the optimal exploitation of the sparsity of the KKT system!
- We also get for free a feedback policy: $\mathbf{u} = -\mathbf{K}\mathbf{x}$!

LQR via Riccati Recursion - Code Example



Thank you

Any Questions?

