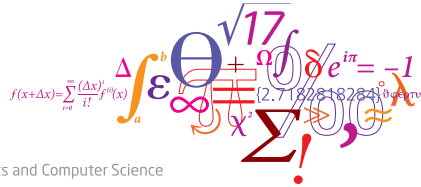


## 02465: Introduction to reinforcement learning and control

Direct methods and control by optimization

Tue Herlau

DTU Compute, Technical University of Denmark (DTU)



DTU Compute  
Department of Applied Mathematics and Computer Science

## Lecture Schedule

### Dynamical programming

- 1 The finite-horizon decision problem  
7 February
- 2 Dynamical Programming  
14 February
- 3 DP reformulations and introduction to Control  
21 February

### Control

- 4 Discretization and PID control  
28 February
- 5 **Direct methods and control by optimization**  
7 March
- 6 Linear-quadratic problems in control  
14 March
- 7 Linearization and iterative LQR  
21 March

### Reinforcement learning

- 8 Exploration and Bandits  
28 March
- 9 Bellmans equations and exact planning  
4 April
- 10 Monte-carlo methods and TD learning  
11 April
- 11 Model-Free Control with tabular and linear methods  
25 April
- 12 Eligibility traces  
2 May
- 13 Deep-Q learning  
9 May

Syllabus: <https://02465material.pages.compute.dtu.dk/02465public>  
Help improve lecture by giving feedback on DTU learn

### Reading material:

- [Her25, Chapter 15]

### Learning Objectives

- Direct methods for optimal control
- Trajectory planning for linear-quadratic problems using optimization
- Trajectory planning using trapezoidal collocation

## Project part 1

- Great job! Part 2 is online
- Survey on course experience on DTU Learn
- A TA caught a minor issue with  $N \rightarrow N - 1$  in the beginning of todays chapter; new version online. Exercise+slides+algorithm not affected.

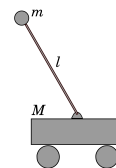
### Recap from last week Dynamics

Dynamics of the form

$$\dot{x}(t) = f(x(t), u(t), t)$$

- $x(t) \in \mathbb{R}^n$  is a complete description of the system at  $t$
- $u(t) \in \mathbb{R}^d$  are the controls applied to the system at  $t$
- The time  $t$  belongs to an interval  $[t_0, t_F]$  of interest

### Recap from last week Example: Cartpole



- Coordinates are  $x = [x \ \dot{x} \ \theta \ \dot{\theta}]$  (angle, angular velocity, cart position, cart velocity)
- Action  $u$  is one-dimensional; the force applied to cart
- Dynamics are

$$\dot{x}(t) = f(x(t), u(t), t)$$

where  $f$  is a fairly complicated function

$$\text{Equality constraint: } x = c \quad (1)$$

$$\text{Inequality constraint: } a \leq x \leq b \quad (2)$$

### Any realistic physical system has constraints

- Simple boundary constraints

$$\mathbf{x}_{\text{low}} \leq \mathbf{x}(t) \leq \mathbf{x}_{\text{upp}}$$

$$\mathbf{u}_{\text{low}} \leq \mathbf{u}(t) \leq \mathbf{u}_{\text{upp}}$$

- End-point constraints:

$$\mathbf{x}_{0, \text{low}} \leq \mathbf{x}(t_0) \leq \mathbf{x}_{0, \text{upp}} \quad (3)$$

$$\mathbf{x}_{F, \text{low}} \leq \mathbf{x}(t_F) \leq \mathbf{x}_{F, \text{upp}}$$

- Time constraints

$$t_{0, \text{low}} \leq t_0 \leq t_{0, \text{upp}} \quad (4)$$

$$t_{F, \text{low}} \leq t_F \leq t_{F, \text{upp}}$$

- The cost function is of the form

$$J_{\mathbf{u}}(\mathbf{x}, t_0, t_F) = \underbrace{c_F(t_0, t_F, \mathbf{x}(t_0), \mathbf{x}(t_F))}_{\text{Mayer Term}} + \underbrace{\int_{t_0}^{t_F} c(\tau, \mathbf{x}(\tau), \mathbf{u}(\tau)) d\tau}_{\text{Lagrange Term}}$$

- Necessary constraint  $-u_{\max} < u(t) < u_{\max}$  and  $\mathbf{x}_0 = [0 \ 0 \ \pi \ 0]$

- Goal is to bring  $\mathbf{x}$  to  $\mathbf{x}^g = [1 \ 0 \ 0 \ 0]$

- Up-right cartpole, version 1:

$$J_{\mathbf{u}}(t_0, t_F, \mathbf{x}) = \|\mathbf{x}(t_F) - \mathbf{x}^g\|^2 + \lambda \int_{t_0}^{t_F} \mathbf{u}(t)^\top \mathbf{u}(t) dt$$

- Constraints  $t_0 = 0, t_F = 3$  (complete in 3 seconds)

- Up-right cartpole, version 2:

$$J_{\mathbf{u}}(t_0, t_F, \mathbf{x}) = t_F - t_0$$

- Constraints  $\mathbf{x}_F = \mathbf{x}^g$

Endless combinations; depends on goal + method you are using

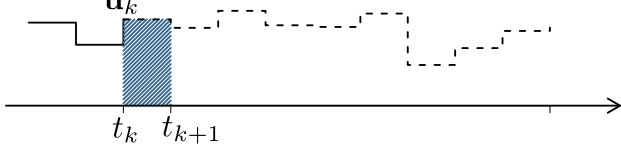
Given system dynamics for a system

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \mathbf{u}(t))$$

Obtain  $\mathbf{u} : [t_0; t_F] \rightarrow \mathbb{R}^m$  as solution to

$$\mathbf{u}^*, \mathbf{x}^*, t_0^*, t_F^* = \arg \min_{\mathbf{x}, \mathbf{u}, t_0, t_F} J_{\mathbf{u}}(\mathbf{x}, \mathbf{u}, t_0, t_F).$$

(Minimization subject to all constraints)



- Simplest choice: Euler's method

- Choose grid size  $N$ :  $t_0, t_1, \dots, t_{k+1} - t_k = \Delta$

- $\mathbf{x}_k = \mathbf{x}(t_k), \mathbf{u}_k = \mathbf{u}(t_k)$

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{f}_k(\mathbf{x}_k, \mathbf{u}_k) \\ &= \mathbf{x}_k + \Delta \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k, t_k) \end{aligned}$$

$$J_{\mathbf{u}=(\mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_{N-1})}(\mathbf{x}_0) = c_f(t_0, \mathbf{x}_0, t_F, \mathbf{x}_F) + \sum_{k=0}^{N-1} c_k(\mathbf{x}_k, \mathbf{u}_k)$$

$$c_k(\mathbf{x}_k, \mathbf{u}_k) = \Delta c(\mathbf{x}_k, \mathbf{u}_k, t_k)$$

- Simple but not very exact

- Last week: Rule-based methods (build  $\mathbf{u}(t) = \pi(\mathbf{x}, t)$  directly)

- Today: Optimization-based methods:

$$\mathbf{u}^* = \arg \min_{\mathbf{u}} J_{\mathbf{u}}(\mathbf{x}_0)$$

- Direct optimization of a discretized version of the problem

- Next week: DP-inspired planning methods

## Infrastructure: Nonlinear program

A non-linear program is an optimization task of the form

$$\begin{aligned} \min_{z \in \mathbb{R}^n} E(z) \quad & \text{subject to} \\ & h(z) = 0 \\ & g(z) \leq 0 \\ & z_{\text{low}} \leq z \leq z_{\text{upp}} \end{aligned}$$

i.e. the objective is to find the  $z$  that minimizes  $E$  under the constraints.

- If problem is not too complex, can use methods such as **sequential convex programming** to find  $z^*$ .
- Requires luck and engineering
  - Needs a good initial guess
  - Improves when given gradient of  $J$  and Jacobian of  $f$  and  $h$ .

## Infrastructure: Linear Quadratic program

A special case of the optimization task:

$$\begin{aligned} \min \frac{1}{2} x^T Q x + c^T x \quad & \text{subject to} \\ & A x \leq b \\ & F x = g \end{aligned}$$

- When  $Q$  is positive definite and the problem is not very large the solution can always be found

## Optimizing the Discrete Problem: Shooting

Consider the simplest form of a discrete control problem

$$x_{k+1} = A_k x_k + B_k u_k + d_k$$

quadratic cost function

$$J_{u_0, \dots, u_{N-1}}(x_0) = x_N^T Q_N x_N + \sum_{k=0}^{N-1} (x_k^T Q_k x_k + u_k^T R_k u_k)$$

- Given  $u_0, \dots, u_{N-1}$ , all the  $x_k$ 's can be found from the system dynamics:

$$x_2 = A_1 x_1 + B_1 u_1 + d_1 = A_1(A_0 x_0 + B_0 u_0 + d_0) + B_1 u_1 + d_1$$

- Problem equivalent to optimizing  $J_{u_0, \dots, u_{N-1}}(x_0)$  (which is quadratic) wrt.  $u_0, \dots, u_{N-1}$
- This method is called **shooting**
- + **A single linear-quadratic optimization problem**
- + **Easy to understand**

## Optimizing the Discrete Problem: Shooting

- General case

$$x_{k+1} = f_k(x_k, u_k)$$

$$J_{u=(u_0, u_1, \dots, u_{N-1})}(x_0) = c_f(t_0, x_0, t_F, x_F) + \sum_{k=0}^{N-1} c_k(x_k, u_k)$$

- Get rid of all the  $x_k$ 's except  $x_0$ :

$$x_2 = f(x_1, u_1) = f(f(x_0, u_0), u_1)$$

So just optimize  $J_{u=(u_0, u_1, \dots, u_{N-1})}(x_0)$  wrt.  $u$

- + **Easy to understand**
- A big, non-linear program (we cannot avoid that for general dynamics)
- - **Unstable: small changes in  $u_0$  can mean big changes in  $x_N$**
- - **Euler's method is imprecise**
- - **No bueno.** To overcome these issues, we have to take a step back

## The continuous-time control problem

Given system dynamics for a system

$$\dot{x}(t) = f(t, x(t), u(t)) \quad (5)$$

Step 1: Must evaluate this ODE somehow

Subject to a number of dynamical and constant path and end-point constraints, obtain  $u : [t_0; t_F] \rightarrow \mathbb{R}^m$  as solution to

Step 2: Computing this integral

$$\min_{t_0, t_F, x(t), u(t)} \underbrace{c_F(t_0, t_F, x(t_0), x(t_F))}_{\text{Mayer Term}} + \underbrace{\int_{t_0}^{t_F} c(x(\tau), u(\tau), \tau) d\tau}_{\text{Lagrange Term}}$$

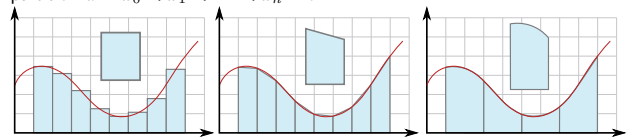
Step 3:  
Minimize over all functions?  
What about constraints?

subject to eq. (5) and whatever constraints are imposed on the system.

**This is a nasty constrained minimization problem**

## Numerical integration

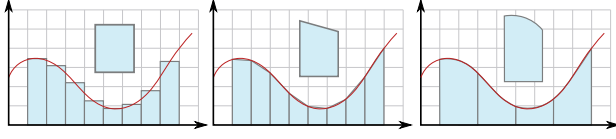
Suppose we wish to approximate a function  $f(x)$ . Divide interval into a partition  $a = x_0 < x_1 < \dots < x_n = b$



Choices corresponds to

- Piecewise constant
- Piecewise linear
- Piecewise 2nd order polynomial (use midpoint to fit the three parameters)

Each provide an approximation for the integral:  $\int_a^b f(x) dx$



- Midpoint rule:  $\approx \sum_{i=0}^{n-1} f\left(\frac{x_{i+1}+x_i}{2}\right) \Delta_i$
- Trapezoid rule:  $\approx \frac{\Delta x}{2} (f(x_0) + 2f(x_1) + 2f(x_2) + \dots + 2f(x_{n-1}) + f(x_n))$
- Simpson's rule:  $\approx \frac{\Delta x}{3} (f(x_0) + 4f(x_1) + 2f(x_2) + 4f(x_3) + 2f(x_4) + \dots + 4f(x_{n-1}) + f(x_n))$

- Given  $t_0$  and  $t_F$  and  $N$
- We discretize the time into  $N$  intervals:

$$t_0 < t_1 < t_2 < \dots < t_{N-1} = t_F$$

- Specifically  $t_k = t_0 + \frac{k}{N-1}(t_F - t_0)$
- For later use we define:

$$h_k = t_{k+1} - t_k, \quad k = 0, \dots, N-2$$

$$\mathbf{x}_k = \mathbf{x}(t_k), \quad k = 0, \dots, N-1$$

$$\mathbf{u}_k = \mathbf{u}(t_k)$$

$$c_k = c(\mathbf{x}_k, \mathbf{u}_k, t_k)$$

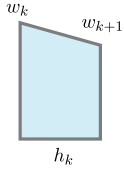
$$\mathbf{f}_k = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k, t_k)$$

**Trapezoid collocation** assumes

$$\int_{t_0}^{t_F} c(\mathbf{x}(\tau), \mathbf{u}(\tau), \tau) d\tau \approx \sum_{k=0}^{N-2} \frac{1}{2} h_k (c_k + c_{k+1})$$

We can at this point evaluate the cost if we know  $\mathbf{x}$  and  $\mathbf{u}$ !

$$c_F(t_0, t_F, \mathbf{x}_0, \mathbf{x}_N) + \frac{1}{2} \sum_{k=0}^{N-2} h_k (c_k + c_{k+1})$$



Recall

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}, \mathbf{u}, t)$$

Integrating both sides

$$\int_{t_k}^{t_{k+1}} \dot{\mathbf{x}}(t) dt = \int_{t_k}^{t_{k+1}} \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t), t) dt$$

Using **trapezoid collocation** we on the right-hand side and integrating the left

$$\mathbf{x}_{k+1} - \mathbf{x}_k \approx \frac{1}{2} h_k (\mathbf{f}_{k+1} + \mathbf{f}_k)$$

- Constraints are translated to simply apply to their knot points:

$$x < 0 \rightarrow x_k < 0$$

$$u < 0 \rightarrow u_k < 0$$

$$\mathbf{h}(t, \mathbf{x}, \mathbf{u}) < 0 \rightarrow \mathbf{h}(t_k, \mathbf{x}_k, \mathbf{u}_k) < 0$$

- Boundary constraints still just apply at boundary:

$$\mathbf{g}(t_0, \mathbf{x}(t_0), \mathbf{u}(t_0)) < 0 \rightarrow \mathbf{g}(t_0, \mathbf{x}_0, \mathbf{u}_0) < 0$$

Optimize over  $\mathbf{z} = (\mathbf{x}_0, \mathbf{u}_0, \dots, \mathbf{u}_{N-1}, t_0, t_F)$

$$\min_{\mathbf{z}} \left[ c_F(t_0, t_F, \mathbf{x}_0, \mathbf{x}_N) + \frac{1}{2} \sum_{k=0}^{N-2} h_k (c_k + c_{k+1}) \right]$$

Such that

$$\mathbf{h}(t_k, \mathbf{x}_k, \mathbf{u}_k) < 0$$

$$\mathbf{g}(t_0, t_F, \mathbf{x}_0, \mathbf{x}_F) \leq 0$$

with convention we iteratively compute  $\mathbf{x}_{k+1}$  from  $\mathbf{x}_k$  starting at  $k = 0$

$$k = 0, \dots, N-2: \quad \mathbf{x}_{k+1} = \mathbf{x}_k + \frac{1}{2} h_k (\mathbf{f}_{k+1} + \mathbf{f}_k)$$

**Wait, did we just solve it?**

## Almost! The final idea:

- Suppose we let  $x_k, u_k$  vary freely (ensure everything can be evaluated)
- But we add the  $N - 1$  constraints:

$$x_{k+1} = x_k + \frac{1}{2}h_k (f_{k+1} + f_k)$$

- The key observation is local changes in  $x_k$  and  $u_k$  have local effects

## Trapezoid collocation method

Optimize over  $z = (x_0, u_0, x_1, u_1, \dots, x_{N-1}, u_{N-1}, t_0, t_F)$

$$\min_z \left[ c_F(t_0, t_F, x_0, x_N) + \frac{1}{2} \sum_{k=0}^{N-2} h_k (c_k + c_{k+1}) \right] \quad (6)$$

$$\text{Such that } z_{lb} \leq z \leq z_{ub} \quad (7)$$

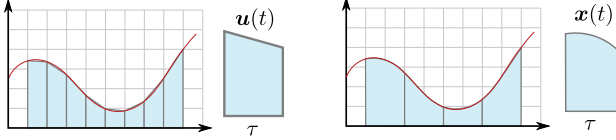
$$h(t_k, x_k, u_k) \leq 0 \quad (8)$$

$$x_k - x_{k+1} + \frac{1}{2}h_k (f_{k+1} + f_k) = 0 \quad (9)$$

- Optimizer also need initial point  $z_0$
- Recall  $f_k = f(x_k, u_k, t_k)$  so last constraint is non-linear

## Reconstruction

Given  $z$ , how do we reconstruct the (predicted) path  $x(t)$  and  $u(t)$ ?



- $u(t)$  was assumed to be linear, using  $\tau = t - t_k$ :

$$u(t) \approx u_k + \frac{\tau}{h_k} (u_{k+1} - u_k)$$

- For  $x(t)$  we assumed

$$\dot{x}(t) \approx f_k + \frac{\tau}{h_k} (f_{k+1} - f_k)$$

- Integrating both sides and using  $x(t_k) = x_k$

$$x(t) = x_k + f_k \tau + \frac{\tau^2}{2h_k} (f_{k+1} - f_k)$$

## Implementation

## Algorithm 1 Direct solver

---

```

1: function DIRECT-SOLVE( $N$ , GUESS= $(t_0^g, t_F^g, x^g, u^g)$ )
2:   Define  $z \leftarrow (x_0, u_0, \dots, x_{N-1}, u_{N-1}, t_0, t_F)$  as all optimization variables
3:   Define grid time points  $t_k = \frac{k}{N-1}(t_F - t_0) + t_0$ ,  $k = 0, \dots, N-1$  ▷ eq. (15.11)
4:   Define  $h_k, f_k = f(x_k, u_k, t_k)$  and  $c_k = c(x_k, u_k, t_k)$ 
5:   Define  $I_{eq}$  and  $I_{ineq}$  as empty lists of inequality/equality constraints
6:   for  $k = 0, \dots, N-2$  do
7:     Append constraint  $x_{k+1} - x_k = \frac{h_k}{2}(f_{k+1} + f_k)$  to  $I_{eq}$  ▷ eq. (15.20)
8:   Add all other path-constraints eq. (15.21) to  $I_{ineq}$  and  $I_{eq}$ 
9:   end for
10:  Add possible end-point constraints on  $x_0, x_F$  and  $t_0, t_F$  to  $I_{eq}$  and  $I_{ineq}$ 
11:  Build optimization target  $E(z) = c_f(t_0, t_F, x_0, x_{N-1}) + \sum_{k=0}^{N-2} \frac{h_k}{2} (c_{k+1} + c_k)$ 
12:  Construct guess time-grid:  $t_k^g \leftarrow \frac{k}{N-1}(t_F^g - t_0^g) + t_0^g$ 
13:  Construct guess states  $z^g \leftarrow (x^g(t_0^g), u^g(t_0^g), \dots, x^g(t_{N-1}^g), u^g(t_{N-1}^g), t_0^g, t_F^g)$ 
14:  Let  $z^*$  be minimum of  $E$  optimized over  $z$  subject to  $I_i$  and  $I_{eq}$  using guess  $z^g$ 
15:  Re-construct  $u^*(t), x^*(t)$  from  $z^*$  using eq. (15.22) and eq. (15.26)
16:  Return  $u^*, x^*$  and  $t_0^*, t_F^*$ 
17: end function

```

---

## Making it work well

- For small  $N$ , method is imprecise, but less sensitive to  $z_0$
- For moderate  $N$ , method is **very** sensitive to  $z_0$
- Initially we do linear interpolation to get  $z_0$
- An idea is to use an optimizer for low value of  $N$ , obtain solution  $z'$
- From this  $z'$ , we can construct  $x'(t)$  and  $u'(t)$
- We run optimizer with higher  $N$  and an initial guess as  $x_k = x'(t_k)$

## Implementation

## Algorithm 2 Iterative direct solver

**Require:** An initial guess  $z_0^g = (x^g, u^g, t_0^g, t_F^g)$  found using simple linear interpolation

**Require:** A sequence of grid sizes  $10 \approx N_0 < N_1 < \dots < N_T$

```

1: for  $t = 0, T$  do
2:    $x^*, u^*, t_0^*, t_F^* \leftarrow \text{DIRECT-SOLVE}(N_t, z_t^g)$ 
3:    $z_{t+1} \leftarrow x^*, u^*, t_0^*, t_F^*$ 
4: end for
5: Return  $u^*, x^*$  and  $t_0^*, t_F^*$ 

```

---

## Implementation:

```

1 # sample.py
2 ineq_cons = {'type': 'ineq',
3             'fun': lambda x: np.array([1 - x[0] - 2 * x[1],
4                                         1 - x[0] ** 2 - x[1],
5                                         1 - x[0] ** 2 + x[1]]),
6             'jac': lambda x: np.array([[ -1.0, -2.0],
7                                         [-2 * x[0], -1.0],
8                                         [-2 * x[0], 1.0]])}
9
10 eq_cons = {'type': 'eq',
11            'fun': lambda x: np.array([2 * x[0] + x[1] - 1]),
12            'jac': lambda x: np.array([2.0, 1.0])}
13
14 from scipy.optimize import Bounds
15 z_lb, z_ub = [0, -0.5], [1.0, 2.0]
16 bounds = Bounds(z_lb, z_ub) # Bounds(z_low, z_up)
17 z0 = np.array([0.5, 0])
18 res = minimize(J_fun, z0, method='SLSQP', jac=J_jac,
19               constraints=[eq_cons, ineq_cons], bounds=bounds)

```

We use sympy because of the gradient/Jacobians

## Example: Pendulum

## Example: Cartpole, the Kelly task

Task is taken from the excellent [Kel17]

- Constraints:  $t_0 = 0, t_F = 2$ , end-point constraints  $x_0$  and  $x_F = x^g$  and  $-20 < u(t) < 20$
- $c(x, u, t) = u(t)^2$
- Grid refinement:  $N = 10$  then  $N = 60$

🔗 lecture\_05\_cartpole\_kelly

## Example: Cartpole, the minimum-time task

From the (also great!) [https://github.com/MatthewPeterKelly/OptimTraj/blob/master/demo/cartPole/MAIN\\_minTime.m](https://github.com/MatthewPeterKelly/OptimTraj/blob/master/demo/cartPole/MAIN_minTime.m)

- Constraints:  $t_0 = 0, t_F > 0$ , end-point constraints  $x_0$  and  $x_F = x^g$  and  $-50 < u(t) < 50$
- $c(x, u, t) = t_F - t_0$
- $N = 8, 16, 32, 70$

🔗 lecture\_05\_cartpole\_time

## Optimizing the Discrete Problem - Collocation

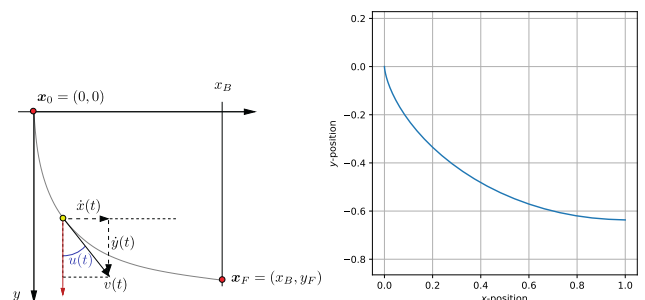
- We can also optimize over both action/state values

The optimisation problem is then defined as

$$\begin{aligned}
 &\text{minimize} && x_N^T Q_N x_N + \sum_{k=0}^{N-1} (x_k^T Q_k x_k + u_k^T R_k u_k) \\
 &\text{subject to} && F'x \leq h' \\
 &&& F''x \leq h'' \\
 &&& A_k x_k + B_k u_k + d_k - x_{k+1} = 0
 \end{aligned}$$

## Example: Brachistochrone

What is the fastest path for a bead to travel  $x_B$  distance in the  $x$ -direction?

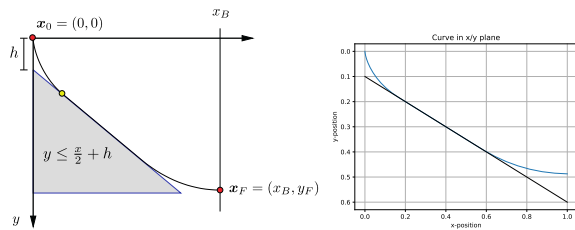


- Cost:  $\min t_F$
- Actions is the angle  $u(t)$ . Dynamics:

$$\dot{x} = v \sin u, \quad \dot{y} = v \cos u, \quad \dot{v} = g \cos u \quad (10)$$

**Example: Brachistochrone with dynamical constraints**

Same as before but bead cannot pass through solid object



- Dynamical constraint

$$h(x) = y - \frac{x}{2} - h \leq 0 \quad (11)$$

**Extra: Hermite-Simpson**

Hermite-Simpson collocation refers to replacing the Trapezoid rule

$$\int_{t_0}^{t_F} c(\tau) d\tau \approx \sum_{k=0}^{N-1} \frac{h_k}{6} (c_k + 4c_{k+\frac{1}{2}} + c_{k+1})$$

For dynamics

$$\mathbf{x}_{k+1} - \mathbf{x}_k = \frac{1}{6} h_k (\mathbf{f}_k + 4\mathbf{f}_{k+\frac{1}{2}} + \mathbf{f}_{k+1})$$

- Generally better for small  $N$
- Scales worse in  $N$



Tue Herlau.

Sequential decision making.

(Freely available online), 2025.



Matthew Kelly.

An introduction to trajectory optimization: How to do your own direct collocation.

*SIAM Review*, 59(4):849–904, 2017.

(See [kelly2017.pdf](#)).