

# Reactive Control

• Differential Kinematics: relationship b/w velocity of all the joints & the velocity of the e.e.  $\Leftarrow$  Jacobian

## • Statics:

$$\text{Power @ joints: } P_z = \boldsymbol{\tau}^T \dot{\boldsymbol{q}}$$

$$\text{power @ ee: } P_e = \mathbf{f}^T \dot{\mathbf{p}} + \mathbf{m}^T \dot{\boldsymbol{w}} \quad \text{moment}$$

$$= [\mathbf{f}^T \quad \mathbf{m}^T] \begin{bmatrix} \dot{\mathbf{p}} \\ \dot{\boldsymbol{w}} \end{bmatrix} = \boldsymbol{w}^T \boldsymbol{v} = \boldsymbol{w}^T \mathbf{J} \dot{\boldsymbol{q}}$$

$$\boldsymbol{\tau}^T \dot{\boldsymbol{q}} = \boldsymbol{w}^T \mathbf{J} \dot{\boldsymbol{q}} \Rightarrow \boldsymbol{\tau}^T = \boldsymbol{w}^T \mathbf{J} \Rightarrow \boldsymbol{\tau} = \mathbf{J}^T \boldsymbol{w}$$

Dynamics: non linear equation

$$M(\boldsymbol{q})\ddot{\boldsymbol{q}} + \underbrace{c(\boldsymbol{q}, \dot{\boldsymbol{q}})}_{h(\boldsymbol{q}, \dot{\boldsymbol{q}})} + g(\boldsymbol{q}) = \boldsymbol{\tau} + \underbrace{\mathbf{J}(\boldsymbol{q})^T \boldsymbol{w}}_{\text{if wrench applied @ ee}}$$

$$h(\boldsymbol{q}, \dot{\boldsymbol{q}})$$

## Joint Motion Control

given a reference joint trajectory  $\Rightarrow$  compute  $\boldsymbol{\tau}(t)$

such that  $\boldsymbol{q}(t) \approx \boldsymbol{q}^{\text{ref}}(t)$

1. Open loop control

3. PD+gravity

2. PID control

4. Inverse Dynamics Control

## Open loop control:

$$\ddot{\gamma}(t) = M(\dot{q}_f^{ref}(t)) \ddot{q}_f(t) + h(\dot{q}_f^{ref}, \dot{q}_f(t))$$

2. PID Control: use tracking error; no knowledge of dynamics

$$\ddot{\gamma}(t) = K_p(\dot{q}_f^{ref}(t) - q_f(t)) + K_d(\dot{q}_f^{ref}(t) - \dot{q}_f(t)) + K_i \int_0^t (\dot{q}_f^{ref}(s) - q_f(s)) ds$$

if  $K_p, K_d, K_i \Rightarrow$  positive definite diagonal matrices

if  $\dot{q}_f^{ref} = 0 \Rightarrow$  guarantee convergence

$$M\ddot{q}_f + C\dot{q}_f + q_f = \ddot{\gamma} = K_p e + K_d \dot{e} + K_i \int_0^t e ds$$

(1) S.S.  $\Rightarrow \ddot{q}_f = \dot{q}_f = 0 \Rightarrow e = K_p e + K_i \int_0^t e ds$  derive  
 $\Rightarrow K_p \dot{e} = -K_i e \Rightarrow \dot{e} = -K_p^{-1} K_i e$

## 3. PD + gravity

Use Knowledge of Model

$$\ddot{\gamma}(t) = K_p(\dot{q}_f^{ref}(t) - q_f(t)) + K_d(\dot{q}_f^{ref}(t) - \dot{q}_f(t)) + g(q_f(t))$$

@ Steady state  $\Rightarrow g = K_p e + g \Rightarrow K_p \dot{e} = 0$

4. Inverse Dynamics control same as open loop

But added the error tracking

$$\begin{cases} \ddot{\gamma} = M\ddot{q}_f^{des} + h \\ \ddot{q}_f^{des} \triangleq \ddot{q}_f^{ref} + K_p(\dot{q}_f^{ref} - \dot{q}_f) + K_d(\dot{q}_f^{ref} - \dot{q}_f) \end{cases}$$

$$M\ddot{q}_f + h = M\ddot{q}_f^{des} + h$$

$$\ddot{q} = \ddot{q}^{\text{ref}} + K_p e + K_d \dot{e} \Rightarrow \ddot{e} = K_p e + K_d \dot{e}$$

second order linear

$$\dot{y} = \begin{bmatrix} 0 & I \\ -K_p & -K_d \end{bmatrix} y \quad y \triangleq \begin{bmatrix} e \\ \dot{e} \end{bmatrix}$$

- need  $M, h \leftarrow$  complete knowledge of model
- stable if  $K_p, K_d > 0$  & diagonal

## Ub • Variant of Inverse Dynamics

$$\tilde{T} = M \ddot{q}^{\text{ref}} + h + K_p e + K_d \dot{e}$$

feedback terms are not pre-multiplied by mass

$$M \ddot{q} + X = M \ddot{q}^{\text{ref}} + h + K_p e + K_d \dot{e}$$

$$\ddot{q} = \ddot{q}^{\text{ref}} + M^{-1} K_p e + M^{-1} K_d \dot{e}$$

$$\dot{y} = \begin{bmatrix} 0 & I \\ -M^{-1} K_p & -M^{-1} K_d \end{bmatrix} y$$

## Cartesian-Space Motion Control

Given a reference for e.e  $x^{\text{ref}}(t)$ , compute  $\tilde{T}(t)$   
such that  $x(t) \approx x^{\text{ref}}(t)$

1. Compute  $q^{\text{ref}}(t)$  s.t.  $x^{\text{ref}}(t) = \text{FG}(q^{\text{ref}}(t))$  forward geometry

2. Operational Space dynamics  $\Rightarrow$  Inverse Dynamics  
 3. Inverse Kinematics + Inverse Dynamics

1. Compute  $q_{ref}(t)$  s.t.  $x_{ref}(t) = FG(q_{ref}(t))$   
 → then apply Joint motion control  
 \* restricts to only one solution  
 \* Don't know  $x_{ref}$  in advance  
 \* going from  $q_{ref} \rightarrow x_{ref}$  requires solving NL equations

2. Operational Space Dynamics [os inverse Dynamics]
- same as inverse dynamics control!
  - \* Take dynamics of joint space & project it onto the ee. space using matrix multiplication
  - \* manipulate the equation to get the relationship bt ee acceleration & joint torques!

$$JM^{-1}[M\ddot{q} + h = \ddot{x}]$$

$$J\ddot{q} + JM^{-1}h = JM^{-1}\ddot{x}$$

$$\ddot{x} - J\ddot{q} + JM^{-1}h = JM^{-1}\ddot{x}$$

$$\cancel{\lambda \ddot{x} + \lambda (JM^{-1}h - J\ddot{q}) = \cancel{\lambda JM^{-1}J^T} f}$$

$$\cancel{\lambda} (q)\ddot{x} + M(\dot{q}, q) = f$$

$$\begin{aligned} x &= FG(q) \\ \dot{x} &= J\dot{q} \\ \ddot{x} &= J\ddot{q} + \ddot{J}\dot{q} \\ \ddot{J}\dot{q} &= \ddot{x} - \ddot{J}\dot{q} \\ \lambda &= (JM^{-1}J^T)^{-1} \\ \ddot{x} &= J^T f \end{aligned}$$

$$f = \Lambda \ddot{x}^{des} + \mu$$

$$\ddot{x}^{des} = \ddot{x}^{ref} + K_p(x^{ref} - x) + K_d(\dot{x}^{ref} - \dot{x})$$

$$\tau = J^T \Lambda \ddot{x}^{des} + J^T \Lambda (JM^{-1}h - j\ddot{q})$$

Simplified version

$$\tau = J^T \Lambda \ddot{x}^{des} + h$$

$$JM^{-1}M\ddot{q} = JM^{-1}J^T \Lambda \ddot{x}^{des}$$

$\ddot{x} - j\ddot{q} = \ddot{x}^{des}$  very small

### 3. Inverse Kinematics + Inverse Dynamics

$$\ddot{x} = J\ddot{q} + J\ddot{q} \Rightarrow \begin{cases} \ddot{q}^d = J^T(\ddot{x} - J\ddot{q}) \\ \tau = M\ddot{q} + h \end{cases}$$

$J^+ = J(JJ^T)^{-1}$   
Moore-Penrose  
Pseudo inverse

$$M\ddot{q} + h = M\ddot{q}^d + h$$

$$\ddot{q} = J^+(\ddot{x}^d - J\ddot{q}) \rightarrow J\ddot{q} = J^+ J^+(\ddot{x}^d - J\ddot{q})$$

$$\rightarrow J\ddot{q} + J\ddot{q} = \ddot{x}^d \rightarrow \ddot{x} = \ddot{x}^d$$

[Postural Task]

Redundancy

$$\dim(x) < \dim(q)$$

→ need to extend your control loop  $\Rightarrow$  adding an extra objective, stabilize the motion of the joint system.

$$JM^{-1}[M\ddot{q} + h = \boldsymbol{\gamma}] \Rightarrow J\ddot{q} + JM^{-1} = JM^{-1}\boldsymbol{\gamma}$$

$$\ddot{x} - J\ddot{q} + JM^{-1} = JM^{-1}\boldsymbol{\gamma} \rightarrow L\ddot{x} + \mu = \underbrace{LJM^{-1}\boldsymbol{\gamma}}_{J^{T\#}}$$

$LJM^{-1} \Rightarrow J^{T\#}$  Transpose pseudo inverse

\* need  $\boldsymbol{\gamma}$  that results in  $\emptyset$  when  $J^{T\#}$  is multiplied

$$L\ddot{x} + \mu = J^{T\#}(\boldsymbol{\gamma} + \boldsymbol{\gamma}_0)$$

$$\boldsymbol{\gamma}_0 = (I - JJ^{T\#})\boldsymbol{\gamma}_1$$

Verify:  $J^{T\#} \boldsymbol{\gamma}_0 = \left( J^{T\#} - J^{T\#} J J^{T\#} \right) \boldsymbol{\gamma}_1 = \emptyset$

$$\boldsymbol{\gamma}_1 = M(K_p(q^{\text{ref}} - q) - K_d(\dot{q}^{\text{ref}} - \dot{q})) + h$$

### Joint space Inverse Dynamics

$$\boldsymbol{\gamma} = J^T L \ddot{x}^d + J^T \mu + (I - J J^{T\#})(M \ddot{q}^d + h)$$

# Interaction Control Regulate contact forces

## 1. Direct force Control

\* Desired force  $\rightarrow$  controller to apply that force

## 2a. Indirect / Impedance Control

\* regulate force by controlling relationship bt motion & force.

Desired Dynamics : Mass-Spring-Damper System

$$M\ddot{x} + B\dot{x} + Kx = f \quad \text{①} \quad \leftarrow \text{in Cartesian space}$$

inertia damper stiffness

$$\ddot{q} = \ddot{x} - \ddot{q}$$

Real Dynamics:

$$M\ddot{q} + h = \tilde{\gamma} + J^T f \quad \leftarrow \text{Joint space} * JM^{-1}$$

$$\ddot{q} + JM^{-1}h = JM^{-1}\tilde{\gamma} + JM^{-1}J^T f \quad * \Lambda = (JM^{-1}J^T)^{-1}$$

$$\ddot{x} + \Lambda(JM^{-1}h - \ddot{q}) = \Lambda JM^{-1}\tilde{\gamma} + f$$

$$\ddot{x} + M = J^T \tilde{\gamma} + f \quad \text{②} \quad \leftarrow \text{Cartesian space}$$

from eq 1:  $\ddot{x} = M^{-1}(f - B\dot{x} - Kx)$  substitute in eq 2

$$\begin{aligned} \ddot{x} + M^{-1}(f - B\dot{x} - Kx) + M &= J^T \tilde{\gamma} + f \\ \ddot{x} + M^{-1}(f - B\dot{x} - Kx) + M &= J^T \cancel{f_m} + f \quad \tilde{\gamma} = J^T f_m \end{aligned}$$

$$\Lambda M^{-1} f + \Lambda M^{-1} (-B\dot{x} - Kx) + \mu = fm + f$$

$$\Lambda M^{-1} (-Bx - Kx) + \mu = fm + (I - \Lambda M^{-1}) f \quad \text{problem:}$$

$$fm = \Lambda M^{-1} (-B\dot{x} - Kx) + \mu - (I - \Lambda M^{-1}) f \quad \begin{matrix} \text{need to} \\ \text{measure} \end{matrix}$$

\* If we choose the desired inertia to be close to the real one then  $M_d = \Lambda$

$$\therefore f_m = -B\dot{x} - Kx + \mu$$

$$\tau = J^T (-B\dot{x} - Kx + \mu)$$

## 2b. Simplified Impedance Control

Replace  $\mu$  w/  $h$

$$\tau = h - J^T (B\dot{x} + Kx)$$

↑ analytical  
↓ numerical

## QP-Based Reactive Control

A numerical optimization problem that take into account the limits of the system

$$M(q)\ddot{q} + h(q, \dot{q}) = \tau \text{ given } \dot{q}(t) \rightarrow \text{find } \tau \text{ st } q(t) \approx \dot{q}(t)$$

$$\tau_d = M\ddot{q}_d + h \Rightarrow \ddot{q}_d = \dot{q}^r + K_p e + K_d \dot{e}$$

→ optimization problem:

$$(\tau^*, \ddot{q}^*) = \underset{\tau, \ddot{q}}{\operatorname{arg\,min}} \| \ddot{q} - \ddot{q}_d \|^2 \quad \begin{matrix} \text{How to solve} \\ \text{this?} \end{matrix}$$

$$\text{s.t. } M\ddot{q} + h = \Sigma$$

## Joint position Limits

- \* same approach to vel limits

$$q(t+\Delta t) = q(t) + \Delta t \dot{q}(t) + \frac{1}{2} \Delta t^2 \ddot{q}(t)$$

$$q_{\min} \leq q + \Delta t \dot{q} + \frac{1}{2} \Delta t^2 \ddot{q} \leq q_{\max}$$

- \* could result in large unfeasible accelerations!

## Convex optimization problems

- \* have only one minimum  $\rightarrow$  if local minimum found, it is the global minimum

### Quadratic Programs (QP)

- \* linear equality/inequality constraints

$$\text{* convex quadratic cost: } \frac{1}{2} x^T H x + g^T x, H \succ 0$$

### Least Squares Programs (LSP) subclass from (QP)

- \* linear equality/inequality constraints ( $Ax \leq b$ )

- \* Cost is 2 norm of affine function  $\|Ax - b\|^2$

## QP in Cartesian Space

given reference e.g. trajectory  $x^*(t)$ , find  $\Sigma(t)$

$$\text{s.t. } x(t) \approx x^*(t)$$

$$\ddot{x} = J \ddot{q} + \ddot{J} \dot{q} \quad \ddot{x}^d \triangleq \ddot{x}^r + K_p(x^r - x) + K_d(\dot{x}^r - \dot{x})$$

$$\begin{cases} \ddot{q}^d = J^T (\ddot{x}^d - \ddot{J} \dot{q}) \\ \ddot{x}^d = M \ddot{q} + h \end{cases}$$

$$J^T = J^T (J^T J)^{-1}$$

$$(\tau^*, \dot{q}^*) = \underset{\tau, \dot{q}}{\operatorname{argmin}} \| J\ddot{q} + \dot{J}\dot{q} - \ddot{x}^d \|^2$$

s.t.  $M\ddot{q} + h = \Sigma$

## Task Models

Generalizes e.e. control

\* describes task w/ function  $e(\cdot)$  to minimize

\*  $\underline{e}$  is error bt reference & real trajectory

$$e(x, u, t) = \gamma(x, u) - \underbrace{\gamma^r(t)}_{\text{trajectory}} \quad x \triangleq (q, \dot{q}) \quad u \in \Sigma$$

cannot directly minimize  $(e)$  but we can its derivative in time (in the future) & because it relies on  $q$  &  $\dot{q}$  which are not decision variables in the LSP.

\* 3 Kinds:

- Affine functions of  $u \triangleq \gamma$

- Non-linear functions of  $x \triangleq (q, \dot{q})$

\* for func. of  $\dot{q} \Rightarrow$  impose 1st derivative

\* for func. of  $q \Rightarrow$  impose 2nd derivative

## 1. Velocity Task function

$$e(x, u, t) = e(q, \dot{q}, t) = \gamma(\dot{q}) - \dot{\gamma}^r(t) \quad \text{decision variable}$$

$$\dot{q} = \frac{\partial \gamma}{\partial q} \quad \frac{\partial \gamma}{\partial \dot{q}}$$

$$\dot{\gamma}^r(t) = T\ddot{q} - \dot{\gamma}^r$$

$$\ddot{q} = \frac{\partial q}{\partial t} - \dot{y}^r = J \dot{e}$$

J

affine linear func

choose  $\dot{e}$  as linear feedback:

$$\dot{e} = -K_e e + \gamma_0$$

$$\ddot{J}\dot{q} - \dot{y}^r = -K_e e \Rightarrow \boxed{\dot{e} + K_e e = 0} \text{ minimize}$$

$$J\ddot{q} - [y^r - K_e e] = 0$$

Ax - b = 0

## 2. Configuration Task function

$$e(x, u, t) = e(q, t) = y(q) - y^r(t)$$

$$\dot{e} = \frac{\partial Y}{\partial q} \frac{\partial q}{\partial t} - \dot{y}^r = \ddot{J}\dot{q} - \dot{y}^r$$

$\ddot{q}$  nota decision variable

$$\ddot{e} = \ddot{J}\dot{q} + \dot{J}\ddot{q} - \ddot{y}^r$$

Choose  $\ddot{e}$  as linear dynamics:

$$\ddot{e} = -K_p e - K_d \dot{e}$$

$$\therefore \ddot{J}\dot{q} + \dot{J}\ddot{q} - \ddot{y}^r = -K_p e - K_d \dot{e}$$

$$\dot{e} + K_p e + K_d \dot{e} = 0 \leftarrow \text{minimize}$$

$$\ddot{J}^T - \underbrace{(\ddot{q}^r - \dot{J}\dot{q} - K_p e - K_d \dot{e})}_{A \times b} = 0$$

## Underactuated systems

- Using Task space inverse dynamics

$$\text{minimize } \frac{1}{2} \|A\dot{z} - b\|^2$$

$$z = (\ddot{q}, \dot{r})$$

s.t.

$$\begin{aligned} & \cancel{\begin{bmatrix} M & -I \\ 0 & I \end{bmatrix} z = -h} \\ & \rightarrow \begin{bmatrix} M & -S^T \end{bmatrix} z = -h \quad \text{selector Matrix} \end{aligned}$$

## Task space Inverse Dynamics Interactions

constrain the motion for rigid contact! define it as a func

$$c(q) = \text{const} \quad \Leftarrow \text{contact points don't move} \quad J = \frac{dc}{dq}$$

$$\dot{J}^T q = 0 \quad \Leftarrow \text{contact points velocities are null}$$

$$\ddot{J}^T q + J^T \ddot{q} = 0 \quad \Leftarrow \text{contact points accelerations are null}$$

Introduce force as a decision variable:

$$\underset{(\ddot{q}, \dot{r}, f)}{\text{argmin}} \|A\dot{z} - b\|^2$$

s.t.

$$\begin{cases} M\ddot{q} + rh = S^T \dot{r} + J^T f \\ \ddot{J}^T q + J^T \ddot{q} = 0 \end{cases}$$

## Friction Cone Constraints

This is a quadratic constraint  $\Leftrightarrow$  can't use LSP

$$\|f_T\|^2 \leq \mu \|f_N\|^2 \quad \mu = \text{friction coefficient}$$

Tangential      Normal

approximate w/ linear pyramidal cone:  
 $B_f \leq 0$

## Multi-Task control

$$g_i(z) \triangleq \|A_i z - b_i\|^2 \quad i = 1 \dots N \text{ tasks}$$

$$\begin{aligned} & \underset{\substack{z=(\dot{q}, \ddot{q}, f) \\ \text{s.t.}}} {\text{minimize}} \sum_{i=1}^N w_i g_i(z) = \|Az - b\|^2 \\ & \quad \begin{bmatrix} M & -J^\top & -S^\top \\ J & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \dot{q} \\ f \\ \ddot{q} \end{bmatrix} = \begin{bmatrix} -h \\ -J\dot{q} \\ -S\ddot{q} \end{bmatrix} \end{aligned}$$

## Inverse Dynamics [Geometry]

\* Numerical I.D.G.  $\Rightarrow$  use computer program to find  $\dot{q}$  by evaluating  $F(\cdot)$  several times  
↳ optimization problem

$$\underset{q}{\text{minimize}} \frac{1}{2} \|p(q) - p_d\|^2 = \frac{1}{2} \|e(q)\|^2 \triangleq C(q)$$

$p_d$  = desired e-e position

$P(q) = \text{Forward geometry function}$

1st order optimality conditions

Find where the gradient of the cost is zero

$$\frac{\partial}{\partial q} C(q) = 0 \Rightarrow \frac{\partial}{\partial q} \left( \frac{1}{2} \bar{e}^T e \right) = \bar{e}^T \frac{\partial e}{\partial q} = \bar{e}^T \underbrace{J}_{J^T} = J^T \bar{e} = 0$$

$r(q) = 0 \Leftarrow \text{start w/ guess } \bar{q}$

$r(q)$

\* Do linear expansion on cost  $r(q)$

$$r(q + \Delta q) = r(\bar{q}) + \boxed{\nabla r(\bar{q})^T \Delta q} = 0 \quad \text{find this}$$

$$\nabla r = \nabla (J^T e) = \underbrace{\nabla J^T}_{\text{ignore}} e + J^T \underbrace{\nabla e}_{\frac{1}{J}} \approx J^T J \quad \begin{matrix} \text{IF } N > 3 \\ \uparrow \\ \text{not invertable} \end{matrix}$$

$$\therefore r(\bar{q}) + J^T J \Delta q = 0 \Rightarrow \Delta q = -(J^T J)^{-1} r(\bar{q})$$

Regularization linear expand  $e$  instead of cost

$$C(q) = \frac{1}{2} \|e\|^2 \approx \frac{1}{2} \|\bar{e} + J \Delta q\|^2 = \frac{1}{2} [\bar{e}^T \bar{e} + \bar{e}^T J \Delta q + (J \Delta q)^T \bar{e} + (J \Delta q)^T J \Delta q]$$

$$= \frac{1}{2} \underbrace{\Delta q^T J^T J \Delta q}_{H} + \underbrace{\bar{e}^T J \Delta q}_{g^T} + \frac{1}{2} \bar{e}^T e \Rightarrow \Delta q = -H^{-1} g$$

minimize  $\|e + J\Delta q\|^2 + \frac{1}{2} \|\Delta q\|^2$  regularization  
 $\Delta q$  Cost

$$\Delta q = (J^T J + \lambda I)^{-1} J^T \bar{e} = -J^+ \bar{e}$$

Pseudo inverse

$\lambda > 0$  regularization weight  $\leq$  very small

\* Issues w/ vanilla gauss newton

1 - can get stuck @ zero slopes

2 - can get stuck in cycles

gradient only gives you local info, taking a large step won't guarantee cost going down  $\rightarrow$  use line search

### I.G. step

• Input:  $q^{(i)}$ ,  $p^d$  current guess output  $q^{(i+1)}$  next guess

1. current error:  $P(q^{(i)}) - p^d$

2. Newton step:  $\Delta q = -(J^T J + \lambda I)^{-1} J^T e$

3. line search:  $q_{try} = q^{(i)} + \alpha \Delta q$ ,  $\alpha = 1$   
 $reduction = \|e\| - \|e(q_{try})\|$

while  $reduction < \gamma \alpha \|e\|$

$\alpha = \beta \alpha$ ;  $q_{try} = q^{(i)} + \alpha \Delta q$ ; reduction

$q^{(i+1)} = q_{try}$

# Optimal Control

- Can be used to do planning  $\Rightarrow$  no need for reference traj =
- Can handle constraints
- Objective clearly specified by the cost function

## optimal Control Problem (OCP)

- takes form of an optimization problem but it is not
  - $x(\cdot), u(\cdot)$  are not vectors  $\Rightarrow$  they are trajectories [infinite dimension objects]

$\hookrightarrow$  the constraints are infinitely many because they all must be valid for all instances in range  $[0, T]$

$$\text{Minimize}_{x(\cdot), u(\cdot)} \int_0^T l(x(t), u(t), t) dt + l_f(x(T))$$

↑ running cost      ↑ terminal cost

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

-  $x(0) = x_0$  initial conditions

-  $g(x(t), u(t), t) \leq 0$  path / actuator constraints

## OC Methods

	opt $\rightarrow$ Disc	Disc $\rightarrow$ opt
global	HJB	Dynamic Programming
local	Indirect Method	Direct Methods

# Dynamic Programming

- DP is a discrete time OCP
- Based on the principle of optimality:
  - the sub parts of an optimal trajectory is also optimal

## DP approach

- Cut the minimization Problem into smaller minimization problems using Bellman's Principle
- $x(-)$  &  $u(-)$  are trajectories in discrete time (matrices)
- Discrete time Dynamics

$$\underset{x, u}{\text{minimize}} \sum_{i=0}^{N-1} l(x_i, u_i)$$

next state is a func of

$$\text{s.t. } x_{i+1} = f(x_i, u_i) \leftarrow \text{current state \& control}$$

→ Split into 2 & move the constraints into the cost of the function  $\Leftrightarrow$  AS AN indicator function

Indicator func  $\begin{cases} 0 & \text{if input=0} \\ \infty & \text{otherwise} \end{cases}$  makes Constraints a very high cost if violated

$$\underset{u, x}{\text{Minimize}} \left[ \sum_{i=0}^{M-1} (l(x_i, u_i) + I(f(x_i, u_i))) + \sum_{i=M}^{N-1} (f(x_i, u_i) + I(f(x_i, u_i))) \right]$$

Indicator func      dynamics constraint

$$\underset{x, u}{\text{Minimize}} \left[ c_0(x_{1:M}, u_{0:M-1}) + V_M(x_M, x_{M+1:N}, u_{M:N-1}) \right]$$

separated because it is a parameter  
not a value

$$\underset{x, u}{\text{Minimize}} c_0(x_{1:M}, u_{0:M-1}) + V_M(x_M)$$

Optimal solution:

$$V_i(x_i) = \underset{u_i}{\text{minimize}} \left[ l(x_i, u_i) + V_{i+1}(f(x_i, u_i)) \right]$$

$V_M(x_M)$   $\Rightarrow$  Value function / optimal cost to go  
↑

optimal cost that we have to pay if we start from a state  $x_M$  onward behaving optimally.

## DP Algorithm

- create the value function for last time step

$$V_N(x_N) = l_f(x_N)$$

- go backwards in time using recursive optimality

$$V_i(z) = \underset{u}{\text{minimize}} \quad l(z, u) + V_{i+1}(f(z, u))$$

minimize only over one control value

- compute optimal control:

$$u_i^*(z) = \underset{u}{\text{argmin}} \quad l(x, u) + V_{i+1}(f(x, u))$$

Q function in RL

not just optimal control trajectory but also an optimal

feed back policy!  $\Rightarrow u$  is a func of  $x$ !

result is a function [not a value]  $\Rightarrow$  Parametric optimiz.

## Numerical Integration

Ordinary Differential equation (ODE)

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

- dependency of  $u$  has been omitted because  $u(t)$  is a function of  $t$  and/or  $x$
- if  $f$  is lipschitz continuous [first derivative of the func. is bounded] there exists a unique solution

### 1. Explicit Euler

simplest numerical integration scheme

$$\dot{x} = \lim_{h \rightarrow 0} \frac{x(t+h) - x(t)}{h}$$

$$h \dot{x} \approx x(t+h) - x(t) \Rightarrow x(t+h) = x(t) + h f(x(t))$$

$\Rightarrow h$  must be sufficiently small  $\Leftarrow$  computationally exp.

## 2. Midpoint

- Take a step to reach middle of time step
  - Compute  $x_1$  then reuse value for the whole step
$$K_1 = f(x_0, t_0), \quad K_2 = f\left(x_0 + \frac{h}{2} K_1, t_0 + \frac{h}{2}\right)$$

$$x_1 = x_0 + h K_2$$
  - Uses 2 function evaluations

## 3. Runge-Kutta (RK)

opras of  $x$

$$x_{n+1} = x_n + \sum_{i=1}^q b_i k_i$$

$$k_i = f(x_n + h \sum_{j=1}^q a_{ij} R_j, t_n + c_i h)$$

$q = \text{order of method}$

if  $a_{ij} = 0 \forall j > i \Rightarrow$  explicit

- consistency order  $p$ ; tells you how quickly the error goes to zero

# Direct methods

- Try to find a soln that is locally optimal
  - Discretize OCP & use NLP to find local optimum solns

-> Discretization steps:

  1. parameterize state & control w/ polynomials
  2. Enforce constraints based on a time grid

- Single shooting (sequential)
- 3. Multiple shooting
- 2. Collocation (Simultaneous)

## 1. Single shooting

• it is about implementing a numerical integration scheme - the state is not a variable of the problem, but you need to compute it from the control [a decision variable]

$\hookrightarrow$  Integrate a differential equation which is the dynamics of the system

- discretize control only  $u(t)$  [piece wise constant]

$$u(t) = g_i \forall t \in [t_i, t_{i+1}]$$

• Compute  $x(t)$  by integrating dynamics

• No need for dynamics constraint

$$\underset{u}{\text{minimize}} \int_0^{t_f} l(x(t_i), u(t_i)) dt + l_f(x(t_f))$$

st.  $g(x(t_i), u(t_i), t_i) \leq 0$  path constraints

• To compute the state & the running cost at the same time  $\Rightarrow$  we can use an augmented state!

$$\bar{x} = (x, c)$$

$$\dot{\bar{x}} = \begin{bmatrix} f(x, u, t) \\ l(x, u) \end{bmatrix}$$

$\hookrightarrow$  integrate this instead of  $x$  &  $c$  separately

Integrate the instants of control

Sensitivities: the derivatives of the integration scheme (Gradient)

why? because we need the gradient of the cost function  $\Rightarrow$  which depends both on the state & control  
However, the state also depends on the control. dependency is decided by the integration scheme

cost func. using Euler

$$c(y) = \int_0^T l(x(t), y(t)) dt + l_f(x(T))$$

$$\approx \sum_{i=0}^{N-1} l(x_i, y_i) \cdot h + l_f(x_N) \leftarrow \text{differentiable}$$

$$\frac{\partial c}{\partial y} = \sum_{i=0}^{N-1} \frac{\partial l(x_i, y_i)}{\partial y} \cdot h + \frac{\partial l_f(x_N)}{\partial y}$$

$$\frac{\partial l}{\partial x_i} \cdot \frac{\partial x_i}{\partial y} + \frac{\partial l}{\partial y}$$
$$\frac{\partial l_f}{\partial x_N} \cdot \frac{\partial x_N}{\partial y}$$

depend on integration scheme

How to compute? Integration Scheme

$$x_{i+1} = \Phi(x_i, y_i) \Rightarrow \frac{\partial x_{i+1}}{\partial y} = \frac{\partial \Phi}{\partial y} \cdot \frac{\partial x_i}{\partial y} + \frac{\partial \Phi}{\partial y}$$

• Initialize  $\frac{\partial x_0}{\partial y} = 0$  because input affect future

• could use one of the Runge Kutta methods to calculate

Explicit Euler:  $\bar{\Phi}(x, u) = x + h f(x, u)$

$$\boxed{\frac{\partial \bar{\Phi}}{\partial x} = I + h \frac{\partial f}{\partial x}}$$

$$\frac{\partial \bar{\Phi}}{\partial u} = h \frac{\partial f}{\partial u}$$

## 2. Collocation

- discretize both  $x(t)$  &  $u(t)$
- enforce dynamics on a time grid

$$x(t) = s_i \quad \forall t \in [t_i, t_{i+1}]$$

$$u(t) = y_i$$

- Replace dynamics  $\dot{x} = f(x, u)$

$$\frac{s_{i+1} - s_i}{t_{i+1} - t_i} - f\left(\frac{s_{i+1} + s_i}{2}, y_i\right) = 0$$

↳ collocation

$\dot{x}$  approximation

$x$  approximation

→ Approximate cost integral

$$\int_{t_i}^{t_{i+1}} l(x(t), \dot{y}(t)) dt \approx \frac{l\left(\frac{s_i + s_{i+1}}{2}, y_i\right)}{2} (t_{i+1} - t_i) \\ \triangleq l_i(s_i, s_{i+1}, y_i)$$

Problem formulation:

$$\underset{s, y}{\text{minimize}} \sum_{i=0}^{N-1} l_i(s_i, s_{i+1}, y_i) + l_f(s_N)$$

$$\text{s.t. } s_0 - x_0 = 0 \quad \leftarrow \text{initial conditions}$$

$$c_i(s_i, y_i, \dot{y}_i) = 0 \quad \leftarrow \text{dynamics}$$

$$g_i(s_i, y_i, \dot{y}_i, t_i) \leq 0 \quad \leftarrow \text{actuator limits}$$

•VIP  $\Rightarrow$  this has a very important feature, which is sparsity! the gradient, the jacobian & the hessian of the problem [derivative of cost & const] are sparse matrices. Because they do not depend on many other variables!  $\Rightarrow$  only the  $i^{\text{th}}$  term

use a solver that exploits sparsity.

# 3. Multiple Shooting (middle ground)

- discretize  $u(t)$
- discretize  $x(t)$  only @ a few points in time [coarse grid]
- Compute intermediate values of  $x(t)$  by integrating dynamics

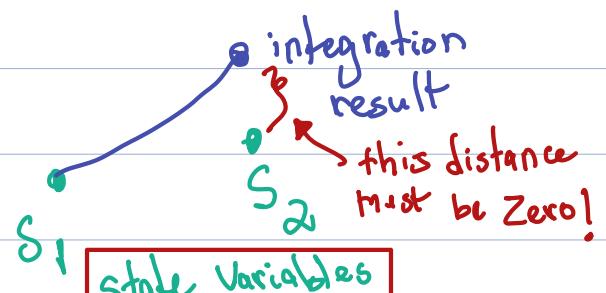
$$\begin{aligned} u(t) &= y_i; \\ \dot{x}(t) &= f(x_i(t), y_i) \end{aligned} \quad \forall t \in [t_i, t_{i+1}]$$

$x_i(t_i) = s_i$  ← few state variables

Numerically compute integral:

$$l_i(s_i, y_i) = \int_{t_i}^{t_{i+1}} l(x_i(t), y_i) dt$$

VIP: Result of the numerical integration might not match w/ the value of the state variable @ the end  
 ↳ continuity constraints must be implemented



→ problem formulation

$$\underset{s, y}{\text{minimize}} \sum_{i=0}^{N-1} l(s_i, y_i) + l_f(x_N)$$

$$\text{S.t. } s_0 - x_0 = 0 \quad g(s_i, y_i) \leq 0$$

continuity constraint  $\Rightarrow$

$$s_{i+1} - \underline{x}_i(t_{i+1})^T s_i y_i] = 0$$

result of integration

- Solver is allowed to violate that constraint while calculating partial solutions.
- Collocation & Multiple shooting find better solutions

## Linear Quadratic Regulator (LQR)

- Method to solve a very specific type of DCP where the dynamics are a discrete time linear system and a quadratic cost function
- Gives an efficient way of solving ocp by exploiting the structure of the problem and solving it as an LSP

$$x_{t+1} = A x_t + B u_t$$

matrices

objectives of LQR is to choose control  $u$  s.t.  $x$  is small &  $u$  is small.

# Competing objectives

Problem formulation:

$$\underset{U}{\text{minimize}} \quad J(U) = \sum_{i=0}^{N-1} \left( X_i^T Q X_i + U_i^T R U_i \right) + X_N^T Q_f X_N$$

$$\text{s.t. } X_{i+1} = A X_i + B U_i$$

quadratic form of

$Q \Rightarrow$  state

$R \Rightarrow$  control

$Q_f \Rightarrow$  final cost

$Q \uparrow \Rightarrow$  small  $X$  more imp.

$R \uparrow \Rightarrow$  small  $U$  more imp.

Rewritten as follows:

$$\underset{(X, U)}{\text{minimize}} \quad J(X, U) = \| \bar{Q} X \|^2 + \| \bar{R} U \|^2$$

Block diagonal mat.

$$\text{s.t. } X = H X_0 + G U$$

$\Rightarrow$  rewrite dynamics as linear eq:

$$X_1 = A X_0 + B U_0 \Rightarrow X_2 = A^2 X_0 + A B U_0 + B U_1$$

$$\begin{bmatrix} X_0 \\ \vdots \\ X_N \end{bmatrix} = \begin{bmatrix} A & & & & \\ & \ddots & & & \\ & & A^2 & & \\ & & & \ddots & \\ & & & & A^{N-1} B \end{bmatrix} X_0 + \begin{bmatrix} B & 0 & \cdots & 0 \\ AB & B & \cdots & \\ \vdots & & \ddots & \\ A^{N-1} B & & & B \end{bmatrix} \begin{bmatrix} U_0 \\ U_1 \\ \vdots \\ U_N \end{bmatrix}$$

a linear func of state & control

• Remove  $X$  from the Problem

$$\underset{U}{\text{minimize}} \quad J(U) = \|\bar{Q}(Hx_0 + GU)\|^2 + \|\bar{R}U\|^2$$

$$= \|\bar{Q}GU + \bar{Q}Hx_0\|^2 + \|\bar{R}U\|^2$$

$$= \|[\begin{matrix} \bar{Q}G \\ \bar{R} \end{matrix}]U + [\bar{Q}Hx_0]_0\|^2$$

decision variable  $\|A X + b\|^2$

LSP: is a square norm of an affine function of a decision variable

LQR w/ DP

• LQR is a rare case where you can use DP because the parametric optimization problem is a quadratic program.

$$V_t(z) = \underset{u_t}{\text{minimize}} \sum_k \left( x_k^T Q x_k + u_k^T R u_k \right) + \left( x_N^T Q_f x_N \right)$$

function of time & state

$\mathbb{1}$  optimal cost starting at state  $z$  @ time  $t$

$V_0(x_0)$  = optimal LQR cost

$$\text{S.t. } x_t = z \quad x_{k+1} = Ax_k + Bu_k$$

→ objective is to recursively minimize the value function starting @ the last time step & going backwards

$$\text{Terminal cost of LQR: } V_N(z) = z^T Q_f z_N$$

→ use bellman principle to go backwards:

$$V_t(z) = \underset{w, u_{t+1}}{\text{minimize}} \sum_i^T \left( Q z_i + w^T R w \right) + \sum_{k=t+1}^{N-1} \left( x_k^T Q x_k + u_k^T R u_k \right) + x_N^T Q_f x_N$$

$$\text{S.t. } x_{k+1} = Ax_k + Bu_k \quad u_t = w \quad \text{next step}$$

$$V_t(z) \underset{w}{\text{minimize}} \sum_i^T Q z_i + w^T R w + V_{t+1}(A z + B w)$$

$$V_{t+1}(z) = z^T P_{t+1} z \quad P_N = Q_f$$

$$V_t(z) = z^T Q z + \min_w (w^T R w + (Az + Bw)^T P_{t+1} (Az + Bw))$$

$$\Rightarrow \nabla f = Qz + 2g = 0$$

Soln =  $w^* = -H^{-1}g$  invertible by definition

~~the def~~

$$R + B^T P_{t+1} B$$

classic soln for  
a QP [eg. IG, Newton]

$$\therefore w^* = K_t z \text{ gain matrix}$$

Summary of LQR using DP:

1. you start by initializing your value function matrix  $P_N$  equal to terminal cost  $Q_f$
2. you go back in time recursively to compute the matrix  $P_{t-1}$  starting from the matrix  $P_t$
3. once you have  $P$  for all the time steps  $\rightarrow$  (value function)  
you can compute the optimal control input for all the time steps by going forward in time  $t=0 \rightarrow t=N-1$   
  - ↳ find the  $K_t$  matrix [gain]

4. once you have  $K_t \rightarrow$  you can compute  $u(t) \rightarrow$  as it is a linear function  $u_t = K_t x_t$   
  - you alternate between using  $u_t$  to find  $x_{t+1}$  & multiplying  $K_{t+1} * x_{t+1}$ ,  $t \rightarrow$  find  $u_{t+1}$

\* we recover an optimal feedback control policy

## Differential Dynamic Programming

- non linear expansion of LQR
- handles non linear Dynamics & cost functions
- approach vaguely resembles Newton methods

### Problem formulation

$$\text{minimize } \sum_{t=0}^{N-1} l(x_t, u_t) + l_N(x_N)$$

$$\text{s.t. } x_{t+1} = f(x_t, u_t) \quad x_0 = x^{\text{init}}$$

1. Start w/ initial guess for  $u$
2. linearize around current trajectory
3. Solve LQR to get variation of  $U$  wrt previous guess
4. line search to ensure convergence

→ a heuristic, no guarantee of convergence

- Apply DP → But we can't apply to the original NL system. Instead we apply to a linearization [local approximation] of the system

$$V(z, u) = \underset{w}{\text{minimize}} [l_i(z, w) + V(f(z, w), i+1)]$$

$$= \underset{w}{\text{minimize}} Q(z, w)$$

• Take a quadratic approximation of  $Q$  function  
 & minimize it.

→ DDP problem is reduced to minimizing the local quadratic model of  $Q$

Soln:  $w^* = \underset{w}{\text{argmin}} Q(z, w) = -Q_w^{-1} (Q_u + Q_{ux} z)$

↗ regularize

- Bias term ↙

- affine func  $\Rightarrow$  not purely linear

→ more complicated than LQR because it is not purely linear func. of state, but affine.  
 Because cost is not purely quadratic [you have the affine terms [linear]]

$$V = \underbrace{DV}_{\text{const}} + \underbrace{V_x z}_{\text{linear}} + \frac{1}{2} z^T \underbrace{V_{xx} z}_{\text{Quadratic}}$$

## Summary of DDP

1. compute state trajectory  $\tilde{x}$  using initial guess for  $u$  going forward

2. Backwards pass  $\Rightarrow$  value func - terminal cost

- then loop backwards in time to calculate:
- Derivatives of Q func.
  - optimal control input  $\omega^*$
  - the Value func.
3. forward pass  $\Rightarrow$  don't have state yet  $\Rightarrow$  can't compute U

$\Rightarrow$  Do a line search starting  $\omega/\alpha = 1$   
 $\Rightarrow$  if cost didn't decrease enough  $\Rightarrow$  take smaller step

Simulate  $\omega$ :  $U = \bar{U} + \omega$  by adding to the previous  
 $= \bar{U} + \alpha \bar{\omega} + K(z)$  guess of  $\bar{U}$   
 $= \bar{U} + \alpha \bar{\omega} + K(x - \bar{x})$

only applied to the feed forward part! not on feedback part.

\* assign  $\bar{x} = x$  &  $\bar{U} = U$

\* obtain optimal open loop control & a locally optimal linear feedback gain K  
 $U = \bar{U} + K(x - \bar{x})$

## Regularization

tune regularization value automatically

$\alpha \uparrow$  slow convergence

- $\mu \downarrow$
- Second derivative  $\bar{Q}_{uu}$  singular
  - $\bar{w}$  too large causing larger control inputs

if  $\bar{Q}_{uu}$  not invertible or  $\frac{J(\bar{U}) - J^*(\bar{U})}{\Delta J(\bar{U})} < c_2$   
 then  $\uparrow \mu$ , otherwise decrease  $\Delta J(\bar{U})$

## Model Predictive Control (MPC)

- use OC for both computing & tracking trajectory
- seen this in DDP  $\rightarrow$  gives feedback gains
- Solve a finite horizon OCP using current state as initial state

$$x^*, u^* = \underset{x, u}{\operatorname{argmin}} \sum_{k=0}^{N-1} l(x_k, u_k)$$

$$\text{s.t. } x_{k+1} = f(x_k, u_k, k) \quad x_0 = x_{\text{meas}}$$

Process:

1. Solve OCP
  2. obtain optimal state & control trajectories
  3. apply first value of control  $u_0$  & throw out the rest
  4. repeat for next iteration
- always behaving optimally at each step

Starting from the last slide:

- Even w/o disturbance, the predicted trajectory using OCP & the one you get from MPC are going to be different!
  - ↳ Because reoptimizing uses a different time horizon in each step

Challenges: • feasibility • stability • computation time

## Infinite Horizon MPC

- All OCPs see the same (infinite) horizon
  - ↳ will all be the same if  $N = \infty$  & cost  $\ell(x, u) \geq \alpha \|x\|$ ,  $\alpha > 0$
- ① Having a finite cost implies stability  
Because it can only be non-infinite if  $x$  goes to zero @ some point
- ② Since predicted & actual are equal  $\Rightarrow$  I have recursive feasibility

- Soln: try to mimic an infinite horizon problem
  - \* add terminal cost & terminal constraint reminiscent of the value function in DP
- Problem formulation:

$$V_N^*(\bar{x}) = \underset{U}{\text{minimize}} \sum_{k=0} l(x_k, u_k) + \underbrace{l_f(x_N)}_{\text{terminal}}$$

st.

$x_{k+1} = f(x_k, u_k)$	$  x_0 = \bar{x}$	cost
$x_k \in \mathcal{X}, u_k \in \mathcal{U}$	$x_N \in \mathcal{X}_F$	terminal constraint

## Feasibility

two cases:

- Input constraints only: always feasible
- Hard state constraints: not just penalties  
if  $N < \infty \Rightarrow$  no guarantee for convergence

Soln #1: [MOAS]  $\leftarrow$  for linear systems

use Maximum output Admissible Set theory

$\hookrightarrow N < \infty$  is enough to enforce recursive feasibility

\* however  $N$  must be sufficiently big

Soln #2: [MCIS]

use Maximum Control Invariant Set

$\hookrightarrow$  can be used as a terminal constraint

• once you are inside the set you can stay inside the set

$\Rightarrow$  hard to compute for non linear systems

Suppose a system w/ linear dynamics

$$\dot{x} = Ax \quad y = Cx$$

- have constraints on output [could be nonlinear]

- objective of algorithm  $\Rightarrow$  find a specific control invariant

set that is also the maximum output admissible set  
↳ this is the set of all initial states  $x$ , s.t. if you start in  $x$   
the constraints on the output are always satisfied

## General Stability theory [Lyapunov Stability]

exponential Lyapunov functions

↳ main tools to prove stability of nonlinear systems

- A set is positive invariant

If  $x$  starts in  $S$  then it stays inside  $S$

↳ Same as control Invariant sets w/o control

\*autonomous dynamics  $\Rightarrow$  no  $u$ .

• if  $S$  is positive invariant  $\Rightarrow$  a value function

$V$  is an exponential Lyapunov func if:

upper & lower bounds:  $V(x) \geq \alpha_1 \|x\|$  ,  $V(x) \leq \alpha_2 \|x\|$

$$V(f(x)) - V(x) \leq -\alpha_3 \|x\|$$

↳ must decrease exponentially along trajectory

- the larger the norm of the state  $\leq$  the more decrease expected

## Stability

- USE Value function [optimal cost to go] as an exponential Lyapunov function

$$V^*(x) = \sum_{i=1}^{n-1} \rho(x_i, y_i) + h_i(x)$$

$$v_0(x_0) = \sum_{i=0}^N l(x_i, u_i) + l_f(x_N)$$

$$\underline{V}_1(x_1) = \sum_{i=1}^N l(x_i, u_i) + l_f(x_{N+1})$$

Consider not optimal  $\Rightarrow$  conservative approach

$\Rightarrow$  difference between the value functions of two consecutive steps:

$$V_1(x_1) - V_0^+(x_0) = l(x_N, u_N) + l_f(x_{N+1}) - l(x_0, u_0) - l_f(x_N)$$

\* if we can prove that this is  $\leq \alpha_3 \|x_0\|$   
then it is an exponential lyapnov func.

## Stability Assumptions

1. origin is @ equilibrium  $\Rightarrow$  running cost = 0  
 $\Rightarrow$  terminal cost = 0

2. union of state & control spaces is a closed set  
 $\Rightarrow$  terminal set is a compact set (closed, bounded)

3. for every state in terminal constraint set  $X_f$ , there must exist a value of  $u$  s.t.  $f(x, u)$  belongs to  $X_f$

\*  $X_f$  is control invariant

\* terminal cost decrease  $l(f(x, u)) - l_f(x) \leq -l(x, u)$

$\hookrightarrow$  tricky to find an  $l_f(x)$  that satisfies this condition

\* running cost must be lower bounded & terminal cost is upper bounded

Note: When we are looking for a local control lyapnov function  $\Rightarrow$  it is only valid inside a certain region.

\* we expect that if we are sufficiently close to the origin, the optimal choice won't be affected by the constraints.

Summary: • Ways to guarantee feasibility.

1. pick  $X_f$  &  $l_f(\cdot)$  s.t.

-  $X_f$  is control invariant at least as much

-  $l_f(f(x,u)) - l_f(x) \leq -l(x,u)$  as running cost

2. can pick  $X_f$  as origin =  $\{0\}$

↳ origin is control invariant

↳ terminal & running costs are zero

} satisfy pt #1

3. can pick  $N$  large enough  $\Rightarrow$  no need for terminal cost or constraint  $\Rightarrow$  Brute force w/ computation power

## Computation Time

• To speed up computation:

- take the optimal soln computed @ time 0 & use it as an initial guess for the problem at next time step.
- instead of iterating until convergence  $\Rightarrow$  just do one optimization  $\Rightarrow$  insures you are doing the computation on the most recent information.

## MPC uncertainties

• there will always be noise. 2 approaches

### Robust MPC

- if you know the bounds of the noise  $\Rightarrow$  apply a safety

Margin wrt the constraint boundary

- easier to implement & a lot more conservative

## Stochastic MPC

- if you know the probability distribution of the noise
  - a lot less conservative

# Reinforcement Learning

- A mathematical framework that allows for solving OCP
- Assume no knowledge of the system's mathematical model
- Tries to find the global optimal policy
- Assumes the dynamics to be stochastic
- Focus on optimization over an infinite horizon

## Markov Decision Process

- Describes the environment in RL problem
- Defined by the tuple  $\langle \mathcal{X}, \mathcal{U}, P, C, \gamma \rangle$
- **Markov property:** the future is independent of the past given the present.
- **State transition Matrix:** a collection of all state transition probabilities  $\nLeftarrow$  irreducible
- **Policy:** Distribution of control inputs over the states. [stochastic]  $\nLeftarrow$  assume deterministic  
 $\pi(u_t | x_t) = P(u_t = u_t | x_{t+1} = x_t)$

Markov Process/Chain: a tuple  $\langle X, P \rangle$   
 $X$ : finite set of states       $P$ : state transition P matrix

- largest eigen value of  $P$  satisfies the following:

$$\min_i \sum_j P_{ij} \leq 1 \leq \max_i \sum_j P_{ij}$$

row

$\Rightarrow$  highest eigen value of the matrix is 1 & all eigen values of  $P$  are smaller than 1

\* state evaluation (Dynamics) uses probability density func  
 $P(x^t | x)$

Markov Reward Process Defined as tuple  $\langle X, P, C, \gamma \rangle$

$C$ : cost func     $\gamma$ : discount factor

$\gamma \rightarrow 0$  myopic eval [future has small impact]

$\gamma \rightarrow 1$  farsight eval [ $\sim$  optimal control]

why?  $\Rightarrow$  avoid inf. cost    uncertainty abt future

Bellman eq:

analytical form:

\*  $V(x) = C(x) + \gamma V(f(x)) \Leftarrow$  for deterministic matrix form:

\*  $V = C + \gamma P V$   $\nabla$  for both deterministic/stochastic  
 - in stochastic systems  $\Rightarrow f(x)$  becomes a probability

distribution & cannot be used to calculate Value func.

## Action Value function (Q)

- **Value func:** a func of the state & yields the ctg starting from that state & following a certain policy -
- **Action value func:** a func of state & control . yields the ctg starting from that state & applying that control & from that point onward follow a certain policy .

$$Q^\pi(x, u) = l(x, u) + \gamma V^\pi(f(x, u))$$

$$V^\pi(x) = Q(x, \pi(x)) \text{ if control from the policy}$$

## optimal value & Q functions

- choose the minimum value over all possible policies !

$$V^*(x) = \min_{\pi} V^\pi(x) , Q^*(x, u) = \min_{\pi} Q^\pi(x, u)$$

- One way to compute optimal Policy:

→ minimize optimal action value func  $Q^*$  over  $u$  !

$$\pi^*(x) = \operatorname{argmin}_z Q^*(x, z)$$

missing info  $\Rightarrow = \operatorname{argmin}_{z \in \mathcal{A}} l(x, u) + \gamma V(f(x, u))$   
Can't be calculated  $\qquad \qquad \qquad$  Dynamics unknown  
+ if we know  $Q^*(x, u) \Rightarrow$  use directly to calculate  $\pi^*$

Bellman optimality eq:

$$V^*(x) = \min_u Q^*(x, u) = \min_u l(x, u) + \gamma V^*(f(x, u))$$

$$\begin{aligned} Q^*(x, u) &= l(x, u) + \gamma V^*(f(x, u)) \\ &= l(x, u) + \gamma \min_{u'} Q^*(f(x, u), u') \end{aligned}$$

→ non linear systems ← minimization is NL

## Bellman operators

- Bellman operator  $T^\pi$ :  $T^\pi(V) = C^\pi + \gamma P^\pi V$

\* input is a value function & output is another value function

- Bellman optimality operator  $T$ :  $T(V) = \min_{u \in U} C^u + \gamma P^u V$

\* Do not have a policy to follow, you minimize over all possible control inputs for all possible states.

→ it looks like the Bellman optimality principle

$$V_K(x) = \min_u l(x, u) + V_{K+1}(f(x, u))$$

## Model Based Control: Dynamic Programming

• in DP, we assume full knowledge of the Model (MDP)

• 2 classes of problems:

1 - prediction & 1 algorithm

→ you have MDP, Policy  $\pi$

→ need to find the value function of  $\pi$

## Iterative Policy Evaluation

→ Apply Iterative Bellman expectation operator

$$V_{K+1} = f(x, \pi(x)) + \gamma V_K(f(x), \pi(x))$$

$$\hookrightarrow V_{K+1} = T^\pi V_K \quad K = \text{iteration \#}$$

• Compute a sequence of approximations of  $V$ , where each approximation is computed from the previous approximation

⇒ process:

1. start w/ arbitrary  $V_0$

2. apply bellman operator iteratively

$$V_{K+1} = T^\pi V_K$$

3. guarantee to converge while  $K \rightarrow \infty$

\* proof & prove that  $V_\pi$  is a unique fixed point of  $T^\pi$

⇒  $V^\pi$  is the only point not affected by  $T^\pi$

1. prove that  $T^\pi$  is contracting:

$$\|T^\pi V - T^\pi Z\|_\infty = \|C^\pi + \gamma P^\pi V - C^\pi - \gamma P^\pi Z\|_\infty$$

$$= \gamma \|P^\pi(V - Z)\|_\infty \quad P^\pi \text{ only changes the order of vectors} \\ \leq \gamma \|V - Z\|_\infty \quad \hookrightarrow \text{no change to max value}$$

2. Starting from  $V_0$ . ⇒ Prove that  $V_K \rightarrow V_\pi$  as  $K \rightarrow \infty$

$$\|V_K - V_\pi\|_\infty = \|T^\pi V_{K-1} - T^\pi V_\pi\|_\infty$$

$$\leq \gamma \|V_{K-1} - V_\pi\|_\infty \dots$$

$$\leq \gamma^K \|V_0 - V_\pi\|_\infty$$

3. Prove  $V_\pi$  is unique fixed point of  $T^\pi$

assume  $T^\pi V = V$ ,  $T^\pi Z = Z$

$$\|T^\pi V - T^\pi Z\|_\infty \leq \gamma \|V - Z\|_\infty$$

$T^\pi$  is contracting

$$\|T^\pi v - T^\pi z\|_\infty = \|v - z\|_\infty$$

Both are fixed

$$\therefore \|v - z\|_\infty \leq \gamma \|v - z\|_\infty$$

- Q soln: 1.  $\gamma > 1$  X impossible because  $\mathcal{T}[0,1]$   
 2.  $v = z$  ✓ only holds when this is true

## 2. Control L \* 2 algorithm

→ you have an MDP

→ need to find the optimal  $v^*$  &  $\pi^*$

**Policy Iteration** ⇒ find optimal policy

process:

1. Start w/ arbitrary  $\pi_0$

→ evaluate policy  
 improve policy ←

2. iterate for  $K=0 \rightarrow \infty$

a. Policy Evaluation  $\pi_K \Rightarrow v^{\pi_K}$

b. improve policy acting greedily wRT  $v^{\pi_K}$

$$\pi_{K+1}(x) = \arg \min_u [l(x, u) + \gamma v^{\pi_K}(f(x, u))]$$

3. as  $K \rightarrow \infty$   $\pi_K = \pi^*$

\* Proof: By showing that each iteration of the algorithm you obtain a policy that is better than the previous one

$$\min_u [l(x, u) + \gamma v^{\pi}(f(x, u))] \leq l(x, \pi_K(x)) + \gamma v^{\pi}(f(x, \pi_K(x)))$$

$\pi_{K+1} \leq \pi_K$   
 one specific control

you have a bigger pool of controls

⇒ you can only be as bad or better

$$\min_u [Q^\pi(x, u)] \leq v^\pi(x, \pi(x))$$

$$Q^\pi(x, \pi') \leq V^\pi(x)$$

## Modified Policy Iteration

- use  $m_K$  policy evaluation iterations to compute  $V^\pi$ 
  - if  $m_K = \infty \rightarrow$  Policy iteration
  - if  $m_K = 1 \rightarrow$  Value iteration

## Value Iteration

- Compute optimal value function
- use to compute optimal policy

### process:

1. start w/ arbitrary estimation of  $V_0$

2. iterate for  $K=0 \rightarrow \infty$

• update value func for all states

$$V_{K+1}(x) = \min_u l(x, u) + \gamma V_K(l(x, u))$$

$$V_{K+1} = T V_K$$

3. guaranteed to converge as  $K \rightarrow \infty$

4. use  $V^*$  to calculate  $\pi^* = \arg \min_u l + \gamma V^*$

## Model free Prediction No control yet

- we have a fixed policy & we just need to find the value function of that policy.
- vaguely resembles DDP approach

## Monte Carlo (MC)

Estimates  $V^\pi$  as the average cost to go

Process:

1. start from state you want to evaluate

2. run the policy

3. collect all costs until the end of episode (must end)

4. sum all costs  $\Rightarrow$  mean estimate of value of state

$\rightarrow$  because stochastic  $\rightarrow$  need to take multiple samples

$\rightarrow$  only applies to episodic MDPs

episode: one execution of Policy from start to terminal state

$$\text{total discounted cost } J_t = l_t + \gamma l_{t+1} + \gamma^2 l_{t+2} + \dots$$

$$\hookrightarrow V^\pi(x) = \mathbb{E}[J_t | x_t = x]$$

in deterministic  $\rightarrow$  1 sample / no need for expected op

first-visit MC:

increment only the first time you see a specific state.

1. counter  $\Rightarrow N(x) \leftarrow N(x) + 1$

2. total cost  $\Rightarrow C(x) \leftarrow C(x) + J(x)$

$$3. V(x) = \frac{C(x)}{N(x)}$$

every visit MC: Counter & total cost are incremented every time a state is visited.

$$\text{* incremental update: } V_N = V_{N-1} + \frac{1}{N} (J_N - V_{N-1})$$

$\rightarrow$  use a fixed value instead of  $N$

$$V(x) \leftarrow V(x) + \alpha (J - V(x))$$

$\alpha \leftarrow$  how fast to update based on most recent sample.

## Temporal Difference Learning - TD<sub>0</sub>

estimates the cost to go by using 1 step look ahead

$$V(x_t) \leftarrow V(x_t) + \alpha_t (\underbrace{r_t + \gamma V(x_{t+1}) - V(x_t)}_{\text{TD error}})$$

TD target  $\Rightarrow$  estimated ctg      TD error

Based on my current estimation of valuefunc

$$\text{MC: } V(x_t) \leftarrow V(x_t) + \alpha_t [J_t - V(x_t)]$$

$$\text{CTG } J_t \approx r_t + \gamma V(x_{t+1}) \quad \text{Bellman eq?}$$

\* Policy evaluation is TD<sub>0</sub> w/  $\alpha = 1$  and sampling all states uniformly

- TD<sub>0</sub> is stochastic approximation

- Since Bellman operator has a unique fixed point  $\rightarrow$  TD<sub>0</sub> converges to  $V^\pi$

- convergence guaranteed if  $\alpha$  is [Robbins Monro]

## MC vs. TD<sub>0</sub>

- MC is unbiased but has high variance
- TD<sub>0</sub> is biased but has low variance

## TDC<sub>></sub>

- use all the step-n cost to go algorithms & add them together using weights.

$$J_t^{\lambda} = (1-\lambda) \sum_{n=0}^{n-1} \lambda^n J_t^{(n)}$$

**Forward View:**  $V(x_t) \leftarrow V(x_t) + \alpha [J_t^{\lambda} - V(x_t)]$

→ can only learn from complete episodes

**Backward View:**

→ need to introduce Eligibility traces:

- measure how often & how recent you have visited state  $x$

-  $e(x) \uparrow$  every time I visit  $x$ , otherwise it decreases

$$\text{TD error } S_t = l_t + \gamma V(x_{t+1}) - V(x_t)$$

$$\therefore V(x) \leftarrow V(x) + \alpha [S_t e_t(x)]$$

## Model Free Control

	value or Action Value	(value or action value) + policy
MDP	$V(x) / Q(x, u)$	$V(x) / Q(x, u) + \pi$
Known	Value Iteration	Policy iteration
Unknown	Direct Methods	Actor-Critic

- if MDP Known : learn value function  $V$
- if MDP un\_Known: learn Action Value function  $Q$

$$\Rightarrow \pi(x) = \arg \min_u l(x, u) + \gamma V(x')$$

need MDP to calculate ↑

Mimize  $Q$  func instead

$$\pi(x) = \arg \min_u Q(x, u)$$

need input  $u$  that minimizes  $Q$

## \* ON Policy learning:

- follow the policy we want to learn
- learn on the job

## \* OFF Policy learning:

- control the system w/ a different policy
- look over someone's shoulder

$\min \leftarrow \text{output}$   
 $\text{arg}\min \leftarrow \text{input}$

## Direct Methods: Q learning

• off policy method

- only learn Q function
- iteratively update an estimate of the Q function
- still use discrete state & control

→ keep track of  $(x, u, l(x, u), x')$  and use them to improve Q based on Temporal Difference (TD)

$$S(Q) = l + \gamma \min_u Q(x', u') - Q(x, u)$$

, need to find lowest value of Q given  $u'$  as input

$$Q_{k+1}(x, u) = Q_k(x, u) + \alpha_k S(Q_k)$$

$$\mathbb{E}[S(Q)] = TQ - Q = 0$$

→ only one Q Function can satisfy this expression →  
a unique fixed point of the bellman operator.

$$\therefore Q = Q^*$$

• needs exploration

## Exploration Strategies

- Greedy  $\rightarrow$  always choose control w/ best cost
- $\epsilon$  Greedy:
  - \* choose random control w/ probability  $\epsilon$
  - \* choose greedy control w/ probability  $1 - \epsilon$
  - \*  $\epsilon$  decreases over time  $\rightarrow 0$ ,  
explore @ beginning, exploit @ end
- Boltzmann exploration:
  - assign prob to each control proportional to its value
  - choose more often controls that leads to lower cost
  - optimism in the face of uncertainty:
    - keep track of uncertainty of value of each control
    - decrease the uncertainty the more you sample that control
    - when you have a lot of uncertainty, you favor that control

## Actor-Critic Methods: Generalized Policy Iteration

- . Extension of Policy iteration to unknown MDP
- . alternate bt [evaluating the policy & updating Q-function]
  - & improving the policy using the value function you estimated
  - ↗ policy evaluation  
↗ policy improvement ↗
- use Model free prediction Methods : MC, TD $\delta$ , TD( $\lambda$ )
- only use a rough eval of policy.
  - $\Rightarrow$  improve Policy based on approximate Value func.
- On Policy

→ The policy used to generate transitions is not the same that is evaluated & improved

↳ the behaviour policy mixes exploration into the target Policy  
\* No guarantee that each iteration is better than the one before it.

## SARSA

Policy evaluation ← critic

- exploit TD to learn Q from on policy samples
- use  $\epsilon$  greedy policy
- update estimate of Q function as follows:

$$S(Q) = l(x, y) + \gamma Q(x', u') - Q(x, u)$$

$$Q_{k+1}(x, u) = Q_k(x, u) + \alpha_k S(Q_k)$$

→ very similar to Q learning → but no minimization of  $Q(x', u')$ !

→ if policy was fixed  $\Rightarrow$  same equation as [TD]

\* as TD, SARSA can diverge in off-policy

Policy improvement ← Actor

Options:

1. act greedily wRT Q ← computed on the fly
2. use  $\epsilon$ -greedy as behaviour Policy
3. Greedy in the limit w/ infinite exploration (GLIE)  
→ alg is allowed not to be greedy in the limit, but as time  $\rightarrow \infty$  the policy becomes greedy

- \* all state-control pairs are visited infinitely many times
- \* Policy converges to a greedy policy
- \* Robbins Monroe Condition for step size

## Summary

	Policy Eval $V(x) \leftarrow I + \gamma V(x')$	Policy Iter. $V(x) \leftarrow I + \gamma V(x)$ $\pi^{(t)} \leftarrow \operatorname{argmin}_u I + \gamma V(x(u))$	Value iteration $V(x) = \min_u I + \gamma V(x(u))$
TD	TD learning $V(x) \leftarrow V(x) + \alpha S$	Sarsa $Q(x, u) = Q(x, u) + \alpha S(Q)$ $\pi(x) = \operatorname{argmin}_u Q(x, u)$	Q-learning $Q(x, u) = Q(x, u) + \alpha S(Q)$

## Value function Approximation

- Discrete space algorithms don't scale well for large systems
- use function approximation techniques instead of tables to represent  $V$  or  $Q$

$$\hat{V}(x, w) \approx V_{\pi}(x) \quad \tilde{Q}(x, u, w) \approx Q_{\pi}(x, u)$$

- $w$  is the weights/parameters of the function
- this is just a parametric function
- ↳ the learning process now aims to learn those parameters s.t. the function approximator

resembles the actual one

→ reduces the size of the vector you're work w/  
→ any differential approximator can be used:

- In RL, FA are usually NNs
  - linear Combinations - polynomials • Fourier basis
- \* as long as it is parametrized & differentiable problems:

1. trying to approximate a moving target
2. The data are not independently normally distributed [sic]

## Gradient Descent

- most common approach
- take a small step in the direction of gradient of function that you want to minimize

$$\min_w J(w) \Rightarrow w = w - \frac{1}{2} \alpha \nabla_w J(w)$$

$\alpha \Rightarrow$  step size

if step is sufficiently small  $\rightarrow$  converge to a local opt.

## Stochastic Gradient Descent (SGD)

- it is stochastic because the quantity we are trying to minimize is the expectation of the square difference bt the real value function & the approximated one.

- it is an expectation because the value of this

depends on the policy  $\pi \Rightarrow$  we want to take the expectation of  $u$  following the distribution of  $\pi$

$$\begin{aligned} J_w &= \mathbb{E}[(V_\pi(x) - \hat{V}(x, w))^2] \\ \Delta w &= -\frac{1}{2} \alpha \nabla_w J = \alpha [(V_\pi - \hat{V}) \nabla_w \hat{V}] \end{aligned}$$

### Process:

- Run a simulation following a policy  $\pi$  that you want to evaluate.
- if we perform the update by sampling the gradient we want to minimize  $\Rightarrow$  then the expected update is equal to the full gradient update that we would get by calculating the gradient of the expectation.

### Incremental Prediction algorithms

- need a target for our value function.

• For MC: compute cost to go  $\bar{J}_t(x)$

$$\Delta w = \alpha (\bar{J}_t - \hat{V}(x_t, w)) \nabla_w \hat{V}(x_t, w)$$

$\Rightarrow$  converges to a local optimum with NL

$\Rightarrow$  in linear approximators  $\Rightarrow$  problem is convex.

For TD0: Target given by TD target:

$$\Delta w = \alpha (l_t + \delta \hat{V}(x_{t+1}, w) - \hat{V}(x_t, w)) \nabla_w \hat{V}(x_t, w)$$

$\Rightarrow$  converges close to global optimum with linear

→ in NL approximators ↵ not guaranteed because of Bias.

## Incremental Control w/ function approximation

- need to switch from value func to  $Q$  because we don't know the dynamics.

→ Do generalized Policy Iteration (GPI)  
w/ an approximate action value func: (Actor-critic)

$$J(w) = \min_w E_{\pi} \left[ (\hat{Q}(x, u, w) - Q_{\pi}(x, u))^2 \right]$$

→ SGD to find local minimum:

- MC: target is  $\text{ctg} \Rightarrow J_t(x_t, u_t)$
- TD 0: target is  $J_t(x_t, u_t) + \gamma \hat{Q}(x_{t+1}, u_{t+1}, w)$
- use  $\epsilon$  greedy to ensure exploration

## Batch learning

- Learn only after having collected a batch of data
- use the batch to perform an update on ws of nn

$$\min_w \sum_{t=1}^T (V_t - \hat{V}(x_t, w))^2$$

- Can be done using experience replay buffer

## Deep Q Network

- an extension of Q learning to continuous state space & discrete control
  - uses experience replay buffer
  - use SGD to optimize w
  - use fixed Q-targets:
    - use one nn for computing the target
    - use another as the main Q and update its ws
    - optimize MSE between the 2 networks
- minimize  $\mathbb{E}_{w} [\hat{y} + \gamma \min_{a'} Q(x', a', \bar{w}) - Q(x, a, w)]$