

UNDERACTUATED ROBOTICS

MIT 6.832 | Spring 2019

LECTURE 10 Trajectory Optimization

Big Picture up till now:

Dynamic Programming

- General: Linear/nonlinear but with additive cost
- Trained Mesh (low dimension),
(e.g. Reinforcement learning)
- Function approximation?
which is hard to put bounds on error

LQR:

- Scales very well
- Linearization means Soln only valid for some RQ

Lyapunov via SDP/SOS

- Scales well (10+ dimensions)
- Can find Soln for all $x \in \mathbb{X}$
- SDP far just used for stability analysis rather than synthesis Q
 - e.g. Found RQ for $\dot{x} = Ax + Bu$ using LQR by using SOS w/ S procedure to find/ maximize region where $V(x) < 0$ choosing $V(x) = x^T S x$
 C cost-to-go from LQR as good guess

- Cons: limited to 'simple' control fns (works for Polynomial / variable transforms e.g. Sines/Cos) but can't capture complex controls that come out of TO
- Cons: Doesn't know a way to use this to synthesize stabilizing control of inv. pendulum

Other Extreme

- Solve for 1 IC to combat high dimensions rather than global control
- tend to grow linearly with dimensions rather than exponentially

- Basic Trajectory Optimization: Formulation

$$\dot{x} = f(x, u) \quad \min_u \int_{t_0}^{t_f} g(x, u) dt$$

s.t. $x(0) = x_0$
Constraints

- Discrete time Case

$$x[n+1] = Ax[n] + Bu[n]$$
$$\min_u \sum_{n=0}^N g(x[n], u[n])$$

s.t. $x[0] = x_0, \dots$

Can't pass directly to solver
since x not defined completely

$$g(x) = \sum_{n=0}^N x^\top Q x + u^\top R u$$

note: w/o constraints same as LQR

(QP) can still add linear constraints

• QPs can run at 100s Hz

• Many variations that remain

Convex (fast w/ global min)

- Eg min time soln for convex opt
Solve w/ constraint

$x[N] = x_g$ & line search

by keep reducing N

→ Transcriptions to numerical implementation:

1) Direct Transcription ↪ solve all vars at once
add x as extra decision variables (slack vars)
as a fn of u

Shooting b/c solve sequentially in time

2) Direct Shooting Method

Solve $x[n]$ as fn of $x_0, u[0], \dots, u[n-1]$

- Still linear constraint

- fewer decision variables but can be less sparse problem

- more difficult to add constraints

- produce uneven gradients (u , high sensitivity)

• Debate which is better for nonlinear systems
Rugg prefers direct transcription

- Back to Continuous time

$$\dot{x} = Ax + Bu$$

$$x[n+1] = x[n] + \Delta t [Ax[n] + Bu[n]]$$

↑ forward Euler

or Runge Kutta ...

- Better than mesh approach from dynamic programming
b/c easier to quantify error

- Sidenote: For unconstrained Linear system \rightarrow LQR
Constrained " \rightarrow MPC

- Choose not to use non-linear dynamic MPC
b/c risk of not finding a soln

• Sweet spot # params vs Numerical integration accuracy

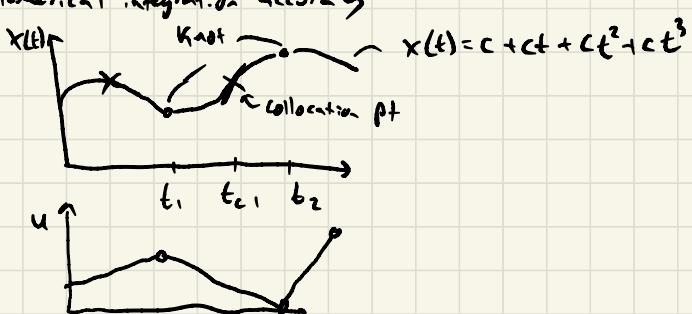
↳ Direct Collocation

$x(t)$ as Cubic Spline

$u(t)$ 1st order hold

- 3rd order accuracy

- Collocation pts - Constr slope
 $\dot{x}(t_c) = f(x, u)$

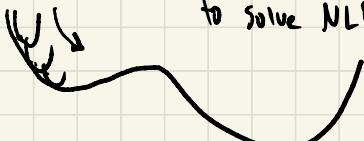


- Note: Can shrink/grow Δt but it creates non-linear optimizer to make them decision vars

- Sidenote: PyDrake Cells auto detects linear vs nonlinear

& calls diff opt algos

- e.g. SNOPT SQP - sequence of QP
to solve NLP



- Sidenote: Can use opt ^{traj} to solve IK of robot end-effectors
- intuition



non-homotopy const \rightarrow
non-convex optimizer

LECTURE 11 Trajectory Stabilization

- Linearize along a trajectory

e.g. pendulum

$$\dot{x} = \begin{bmatrix} \dot{\theta} \\ \frac{1}{m\ell^2} [u - b\dot{\theta} - mg\sin\theta] \end{bmatrix} \xrightarrow[b=0]{g=10} \begin{bmatrix} \dot{\theta} \\ u - 10\sin\theta \end{bmatrix}$$

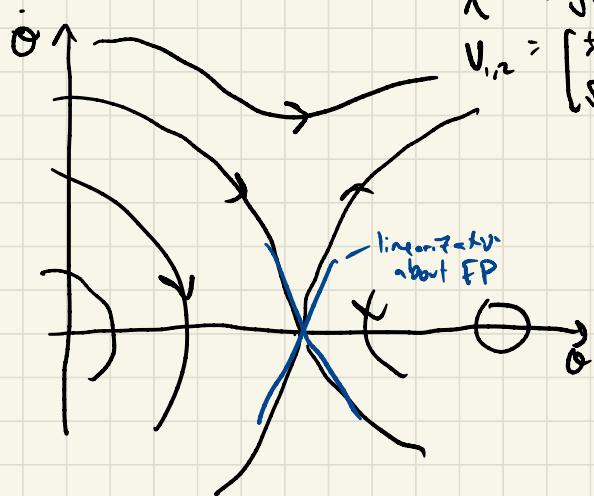
$$\dot{x} \approx f(x_0, u_0) + \frac{\partial f}{\partial x}(x-x_0) + \frac{\partial f}{\partial u}(u-u_0)$$

$$\begin{bmatrix} \dot{\theta}_0 \\ u_0 - 10\sin\theta_0 \end{bmatrix} \xrightarrow[k_{uu}]{k_{xu}} \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 1 \\ -10\cos\theta_0 & 0 \end{bmatrix}$$

$$\lambda = \pm \sqrt{-10\cos\theta_0}$$

$$U_{1,2} = \begin{bmatrix} \pm 1 \\ \sqrt{-10\cos\theta_0} \end{bmatrix}$$



← linearization only good @ fixed
pts in this coord system
(given smooth system)

- Moving Coored System

$x_0(t)$ $u_0(t)$ (e.g. from traj opt)

consistent w/ dynamics

$$\hookrightarrow \bar{x} = x - x_0(t) \quad \bar{u} = u - u_0(t)$$

$$\therefore \dot{\bar{x}} \approx f(\bar{x}, \bar{u}) + \dots$$

$$\dot{\bar{x}} = \underbrace{\frac{\partial f}{\partial x} | \bar{x}}_{A(t)} + \underbrace{\frac{\partial f}{\partial u} | \bar{u}}_{B(t)} \quad \dot{\bar{x}} = A(t)\bar{x} + B(t)\bar{u}$$

↑ powerful covers a lot of control. Eg

$$\bar{x} \rightarrow 0 \text{ then } x \rightarrow x_0(t) \quad \text{LQR works}$$

- Time Varying LQR

t_f ↪ may not be so if from TO output

$$\min_{\bar{u}} \int_0^{t_f} \bar{x}^T Q \bar{x} + \bar{u}^T R \bar{u} dt + \bar{x}(t_f)^T Q_f \bar{x}(t_f)$$

can also make these $f(t)$

$$J(\bar{x}, t) = \bar{x}^T S(t) \bar{x} \quad S(t) > 0 \quad \downarrow \quad \bar{x}^T \dot{S}(t) \bar{x}$$

$$0 = \min_u \left[g(\bar{x}, \bar{u}) + \frac{dJ}{d\bar{x}} f(\bar{x}, \bar{u}) + \frac{dJ}{dt} \right]$$

$$\bar{u} = -R^{-1} B(t) S(t) \bar{x} \quad \text{similar to before}$$

$$= -K(t) \bar{x}$$

↪ boundary value problem defined at t_f so solve it backwards

$$-\dot{S}(t) = S(t) A(t) + A^T S(t) - S(t) B(t) R^{-1} B(t) S(t) + Q$$

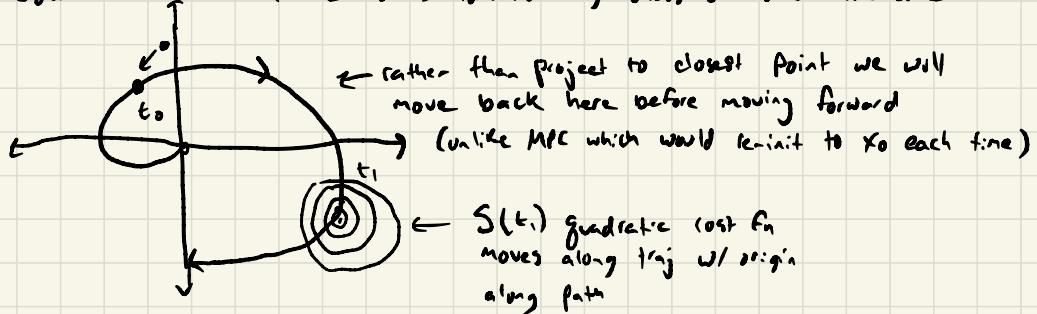
↑ Differential Riccati Equation

$$1) S(t_f) = Q_f$$

$$2) S(t_f) = S \quad \text{final cost-to-go once reach fixed point is cost-to-go of } \infty \text{ t LQR}$$

↪ LQR about vertical fixed pt
smart way to choose final cost

- Downside: ref frame moves forward regardless of where we are



- Unconstrained Linear \rightarrow LQR

- Linear Constraints \rightarrow Linear MPC

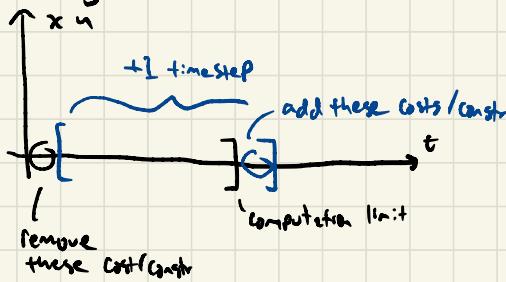
- TV LQR - Locally stabilizes x_0, u_0

- L MPC - locally stabilize Constrained System

- Can't solve QP for $N \rightarrow \infty$

• need finite horizon

- Receding horizon Control



- Guard against infeasibility

- Cost @ end < Cost removed in beginning \Leftarrow Can get this by adding

$$\begin{aligned} &\text{Lyapunov} \\ &\hookrightarrow V(x, k) = \min_u \sum_{\ell=k}^{k+N} x^\top Q x \quad \text{st } \dots \end{aligned}$$

terminal constraint $x[N] = x_G$

guarantees Stability!

$$\hookrightarrow V(x, k+1) \leq V(x, k)$$

- Iterative LQR (for trajectory optimization)

◦ Initial guess $x_0(t), u_0(t)$

but using quadratic approx of nonlinear cost

$$g(x, u) \approx g(x_0, u_0) + \frac{\partial g}{\partial x}(x - x_0) + \frac{\partial g}{\partial u}(u - u_0)$$

$$+ (x - x_0)^T \frac{\partial^2 g}{\partial x^2}(x - x_0) + \dots \quad \leftarrow \begin{array}{l} \text{second order} \\ \text{Taylor} \end{array}$$

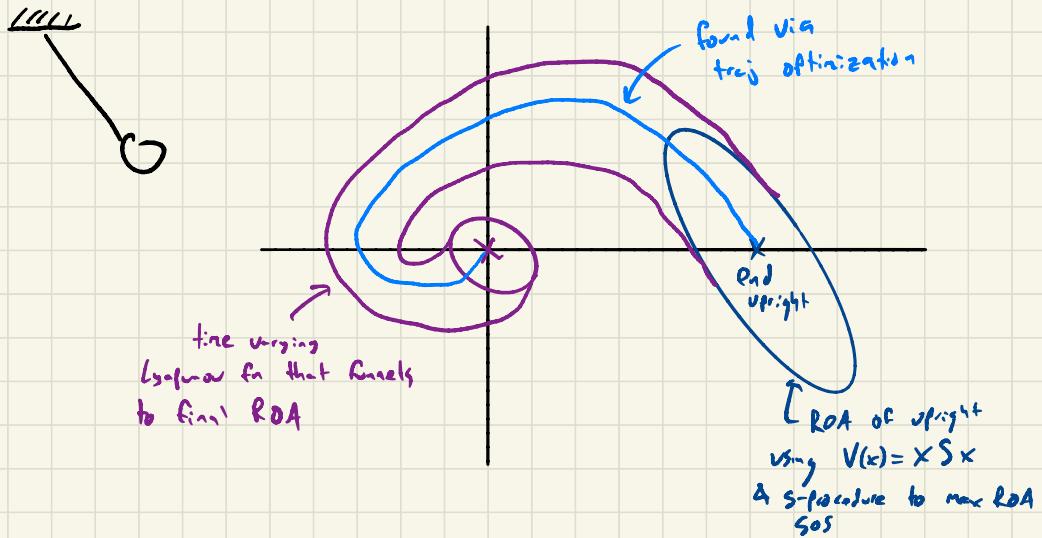
$$\min_u \int_0^T x^T Q x + q_1^T x + q_0 + u^T R u + r_1^T u + r_0 + u^T P x \, dt$$

↳ Intuition: run offline to get trajectory

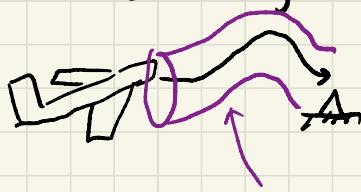
- Start w/ traj & iteratively run ILQR to find new traj to converge to nonlinear cost fn
- quadratic Cost-to-go no longer @ origin
- iteratively find cost fn & linearize to get new cost fn
- like SQP in SNOPT but more clever b/c using

↳ Intuition: MPC - if you can run optimization fast enough you can stabilize
ILQR - if you can stabilize non-trivial trajectory over & over, you can optimize

- Use case of TVLQR & SOS Lyapunov



- Use Case 2 Perching Plane

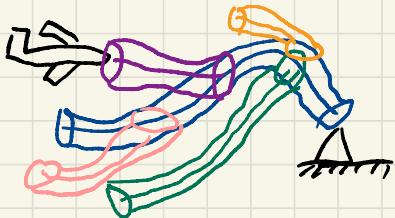


ROA vs SOS
Downside: requires potentially small variation in init condition

Steps:

- System ID w/ motion capture to get nonlinear system dynamics
- Run offline through traj optimization
- Use TVLQR for online feedback stabilization
↳ global control so can run fast (rel MPC)

↳ LQR Trees: combats issue of requiring small initial set to stay in ROA funnel



↳ Fill space with funnels to converge to end ROA

↳ This allows linear fast feedback control (MPC not solution to everything)

↳ Controllability analysis → loose control authority near end → adds flappers

LECTURE 12: Simple Models of Walking

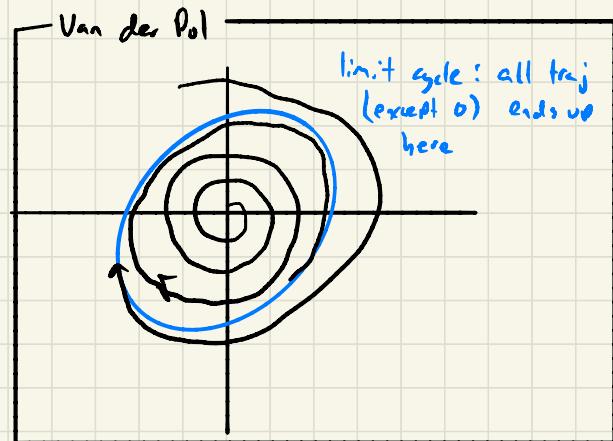
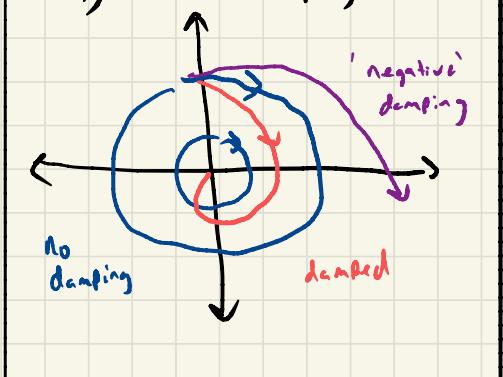
Key Ideas:

- Contact Mechanics - non-smooth mechanics
- Limit cycle stability

- Example (w/o input) Van der Pol oscillator

$$\ddot{\theta} + N(\dot{\theta}^2 - 1)\dot{\theta} + \theta = 0 \quad N > 0$$

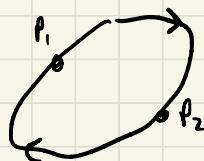
Analogy to linear spring mass:



- Stability of a Trajectory

$$\|x(t) - x^*(t)\| \rightarrow 0 \quad \forall x(t)$$

↪ False for Van der Pol b/c

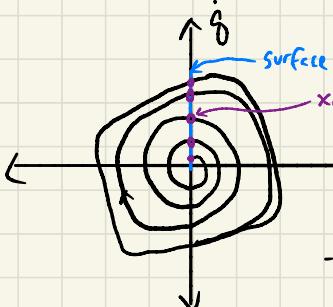


P_1 & P_2 will never converge to same $x^*(t)$

- Limit Cycle Stability

- A periodic trajectory $\forall t=0 \quad x^*(t) = x^*(t+T)$
is limit cycle stable iff (orbital stability)
 $\min_{\tau} \|x(t) - x^*(\tau)\| \rightarrow 0$
↳ find closest point
indep of time

• Surface of Section:



- for more complex systems will require numerical analysis

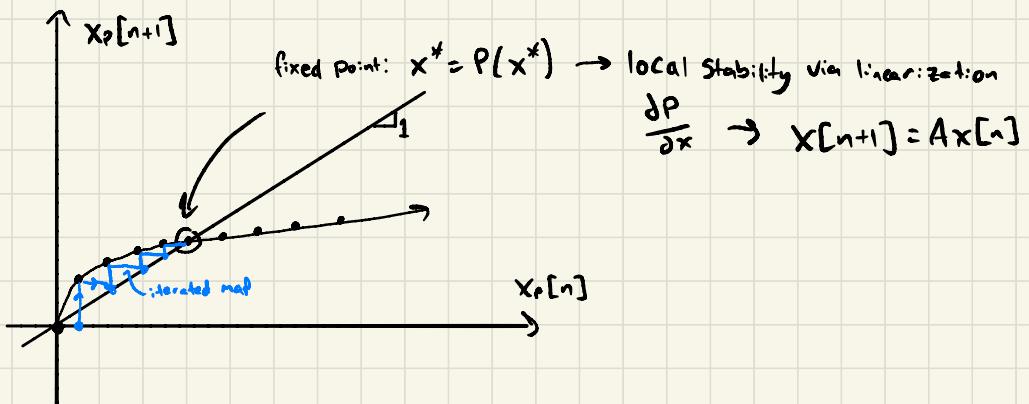
- Section is transverse to the flow
- all trajectories return to the section

$$x_p[n+1] = P(x_p[n])$$

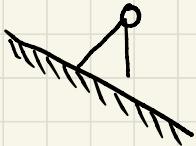
$\uparrow \subset \mathbb{R}^{n-1}$

Poincare Map: transforms limit cycle stability to classical analysis

• Poincare Map for Van der Pol oscillator

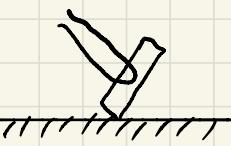


- Contact Models



Rigid vs Soft

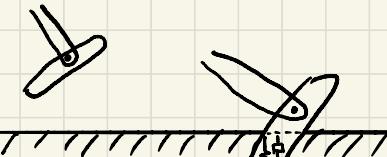
Strict non-penetration



Collisions modelled as impulse force

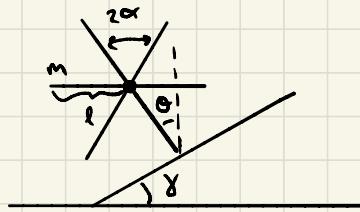
Soft

E.g. Spring damper model



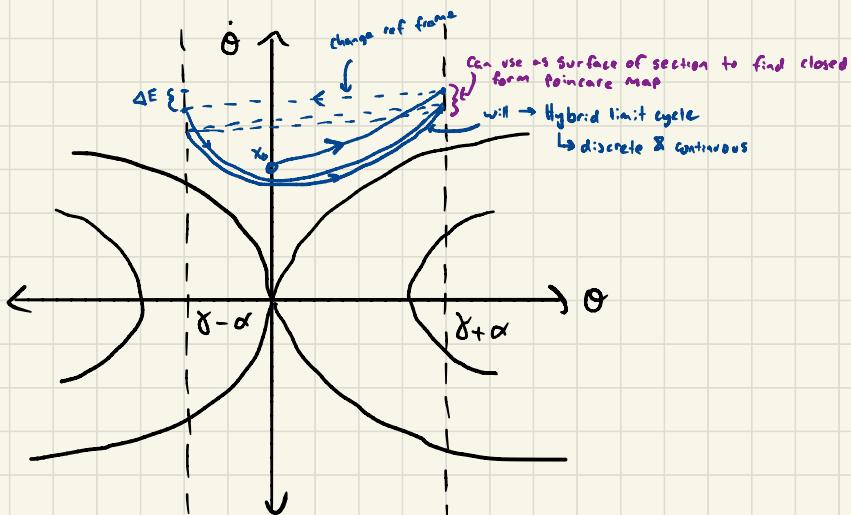
$$\lambda_{\text{normal}} = -k(\text{penetration dist}) - \text{damp}$$

- Rimless Wheel

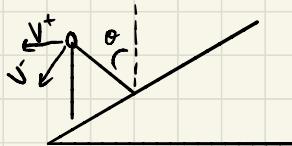


Assumptions:

- No slip (feet on ground as pin joint)
- Collisions are inelastic & impulsive (rigid)
 - ↳ no bouncing
- No double support



- Inelastic Collision



- All energy into ground is lost
- Angular momentum around point of collision is conserved

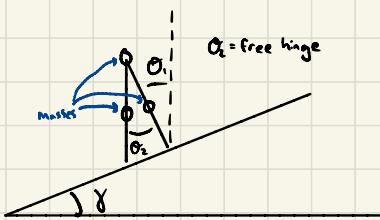
$$L^- = l \times m V^- = m l^2 \dot{\theta}^- \cos(2\alpha)$$

$$L^+ = l \times m V^+ = m l^2 \dot{\theta}^+$$

$$\therefore \dot{\theta}^+ = \dot{\theta}^- \cos(2\alpha)$$

- Side note: Ruvina made this more efficient by having rear foot push off before impact \rightarrow less $\Delta \vec{V}$ after collision (4x more efficient)

- Compass Gait



- Also has stable limit cycle but smaller

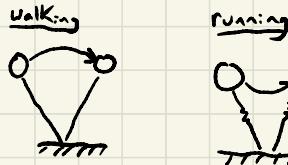
LECTURE 13 Simple Models of Running

- Counterintuitively running (hopping) can be easier than walking

- Definition of Running:

- Existence of an aerial phase

• Exchange of energy e.g.:

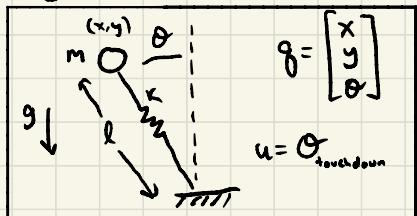


↳ Spring Loaded Inverted Pendulum (SLIP)

- Why Simple Models:

- Tractable
- Mechanical insights
- Comparative Biology
 - ↳ Invariants in data \rightarrow fundamental principles
- As a template for higher DoF robots

- Spring Loaded Inverted Pendulum



Assumptions:

- Massless leg
- ↳ Command θ instantaneously
- Perfectly elastic collision
 - ↳ energy conserved
(threat to stability)
- When foot on ground, pin joint
(∞ sliding friction)

-Poincaré Analysis

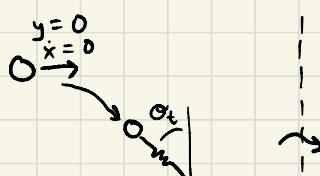
• Apex-to-Apex Map

$$x = \begin{bmatrix} x \\ y \\ \theta \\ \dot{x} \\ \dot{y} \\ \dot{\theta} \end{bmatrix} \leftarrow \begin{array}{l} \text{Given Const } E \quad \dot{x} = f_1(y) \\ = 0 \text{ at apex} \end{array}$$

$$y[n+1] = P_{u[n]}(y[n]) \quad u[n] = \theta_{\text{touchdown}}$$

• Flight Phase

$$x = \begin{bmatrix} x \\ y \end{bmatrix} \quad \ddot{x} = \begin{bmatrix} 0 \\ -g \end{bmatrix}$$



$$y(t_{\text{touchdown}}) = l \cos \theta_{\text{touchdown}}$$

• Standing Phase

$$x = \begin{bmatrix} r \\ \theta \\ \dot{r} \\ \dot{\theta} \end{bmatrix}$$

↳ "minimal coordinates"
reduce Constr. Violation

• Take-off to aerial

$$x(0)$$

$$y(0) \rightarrow y(t_{\text{takeoff}}) = 0$$

$$\dot{x}(0)$$

$$\dot{y}(0)$$

+ Energy Correction at apex (We linearized contact dynamics lose Const E property)

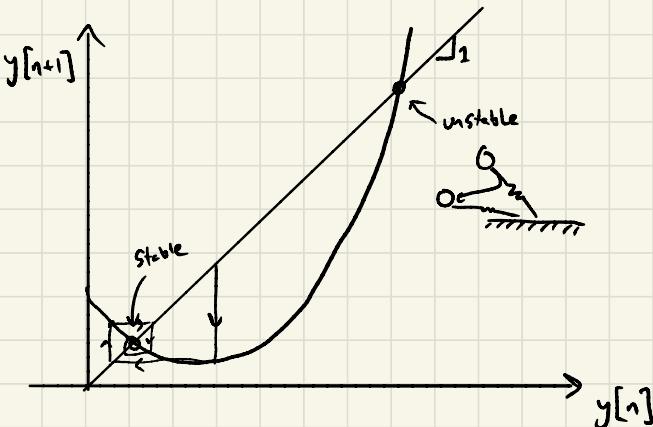
$$mr^2\ddot{\theta} - mr^2\dot{\theta}^2 + mg \cos \theta - K(r_0 - r) = 0$$

$$mr^2\ddot{\theta} + 2mr\dot{r}\dot{\theta} - mgr \sin \theta = 0$$

↳ Hartmut Geyer gave us linearized about small angles

↳ No closed form solution

- SLIP approximate apex-to-apex map:



- SLIP Control

- choose θ touchdown during aerial phase

$$y[n+1] = P(y[n], u[n])$$

- Goal design $u[n] = \pi(y[n])$ to stabilize y^d

- Idea 1: find u^* s.t. $y^d = P(y^d, u^*)$

Linearize P around (y^d, u^*) & do discrete time LQR

- Idea 2: Deadbeat Control (achieve desired y^d in 1 timestep)

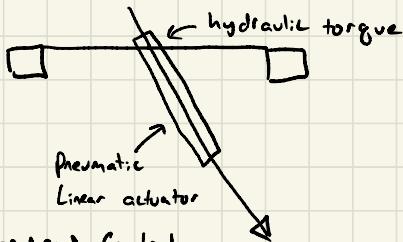
if P is invertible:

if it exists

$$u[n] = P^{-1}(y^d, y[n])$$

↳ Can always set θ to this angle as fn y so it doesn't matter where the ground is. Cool!

- Continuous Control (Raibert Hoppers)



Decomposed Control

1) Hopping height (push at toe-off)

2) Foot touchdown to regulate speed

3) Stabilize attitude during stance

↳ early version of Boston Dynamics's Big Dog used this

↳ now a lot more optimization

LECTURE 14 Computational Methods For Legged Robots

- Simple models:

- Rimless wheel (closed form soln Poincaré map, RQA, limit cycles etc)
- Compass gait (requires numerical tools to understand this & more complex)
- SLLP (running) (requires small & approx to get Poincaré map)

- New concepts

- Limit cycle (stability)
- Hybrid dynamics of contact

- Smooth Systems

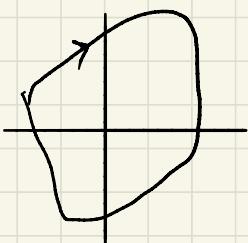
$$\dot{x} = f(x) \quad \text{fixed pt is } f(x^*) = 0$$

↳ may be hard to find but can use optimization

$$\text{find } s.t. f(x^*) = 0$$

• Periodic soln → can also use trajectory optimization

- Ex Van der Pol oscillator (non hybrid case)



$$\text{find } s.t. \dot{x}(t) = f(x)$$

x, t

e.g. direct collocation

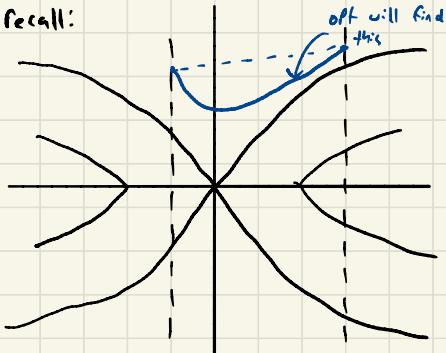
knots in spline
are decision vars
(cycle duration unknown)

$$x(0) = x(t_f)$$

$$x(0) = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \text{to avoid fp @ 0}$$

- Ex Rimless wheel cycle w/ direct collocation (only 1 soln exists for this simple hybrid system)

Recall:



$$m\ddot{\theta} - mg\sin\theta = 0 \quad \dot{\theta}^+ = \cos(2\alpha) \dot{\theta}^-$$

↑ from impact

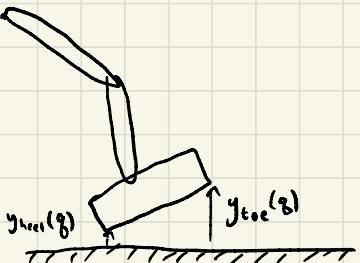
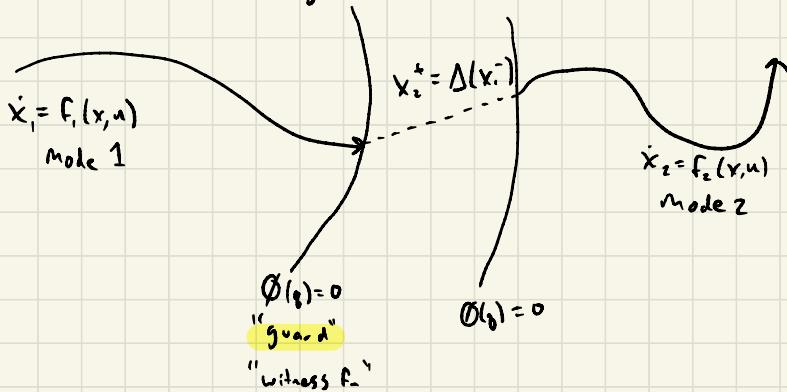
$$\text{find } s.t. \dot{x} = f(x) \text{ via direct}$$

$$\dot{\theta}(0) = \gamma - \alpha$$

$$\dot{\theta}(t_f) = \gamma + \alpha$$

$$\dot{\theta}(0) = \cos(2\alpha) \dot{\theta}(t_f)$$

- More General (**autonomous hybrid Systems**)



- Mode 1 : foot in air
- Mode 2 : heel on ground
- Mode 3 : heel & toe on ground
- Mode 4 : toe on ground

- If you prescribe "**mode schedule**"

$$x[0], \dots, x[k] \quad x[n+1] = f_1(x[n], u[n])$$

$$x[k+1], \dots, x[N] \quad x[n+1] = f_2(x_n, u_n)$$

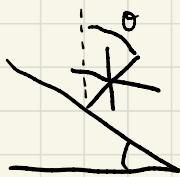
$$\forall k \quad \emptyset(x[k]) > 0$$

$$\emptyset(x(k)) = 0$$

$$x[k+1] = \Delta(x[k])$$

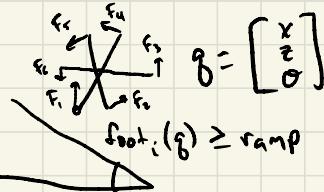
- This method breaks down for more complex systems when order of contacts are unknown

- "minimal coordinates"



vs

"maximal coordinates"



$$M(q)\ddot{q} + C(q, \dot{q})\dot{q} = \tilde{L}_g(q) + Bu + \sum_i J_i(q)F_i$$

$$\frac{\partial \phi_i}{\partial q}$$

contact
jacobian
"contact
forces"

- | if $\phi_i(q) > 0, F_i = 0$
- | if $\phi_i(q) = 0$
- | $\dot{\phi}_i(q) = 0$
- | $\ddot{\phi}_i(q) = 0 \Rightarrow$ solve for F_i to impose the constraint

| ↳ { } time consuming to solve for every combination
can solve w/ opt instead

$$\begin{aligned} \text{find } x \text{ s.t. } \dot{x} = f(x, u, F) \\ \phi(q) = 0 \rightarrow \dot{\phi}(q) = 0 \end{aligned}$$

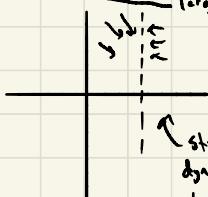
- To avoid mode scheduling for complex systems

◦ Idea 1: Soft Contact

$$F_i = \text{soft-spring-model } (x)$$

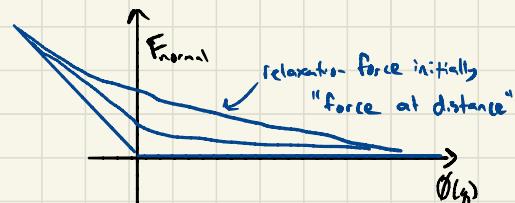
$$= -K(\phi(x)) - \text{damping} \dots$$

K large to be accurate



◦ Stiff dynamics \Rightarrow finding best soln hard
↳ small dt & many knot pts to solve

→ Emo Todorov uses this with relaxation to get soln easier



↳ issue w/ this is too much deformation
could confuse collision engine

• Idea 2: Time Stepping approximation (Big Game engines w/ many objects)

$$x[n+1] = f(x[n], u[n], \lambda[n])$$

↑ backwards Euler update ↓ contact forces

$$x[n+1] = x[n] + dt f(x, u) \leftarrow \text{forward euler}$$

$$\emptyset(g) \geq 0 \quad \lambda[n] \geq 0$$

$$\leq \epsilon$$

Complementarity constraint $\emptyset(g[n]) \cdot \lambda[n] = 0$

↳ same as $\emptyset(g) \geq 0$ or $\lambda = 0$

↳ Surprisingly Effective Solns via Linear Complementarity Problem LCP

↳ "Let the Solve figure out which contact forces are active w/o

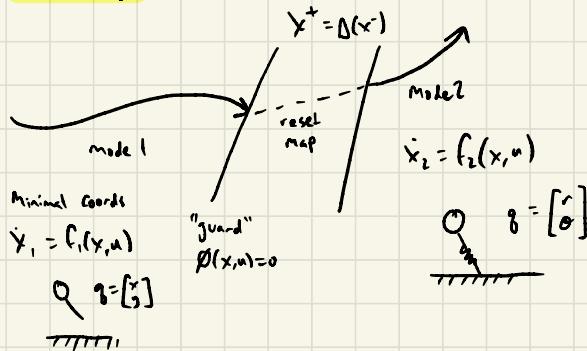
soft contact"

↳ Russ uses this

Lecture 15 Feedback & Lyapunov Through Contact

Last time:

Minimal Coord:



Maximal Coordinates: (hybrid case)

Mode 1

$$\dot{x} = f(x, u, 1)$$

Contact constraint forces

$$g_1(x, u, 1) = 0$$

$$\text{Q } \lambda = 0$$

Mode 2

$$\dot{x} = f(x, u, 2)$$

$$g_2(x, u, 2) = 0$$

$$\text{Q } \lambda = 0$$

Hard Contact

foot(g) = Constant location

$$\frac{d\text{foot}}{dt}(g, \dot{g}) = 0$$

$$g(x, u, 1) \Rightarrow \frac{d^2\text{foot}}{dt^2}(g, \ddot{g}, u, 1) = 0$$

if Soft Contact

$$\lambda = -K \text{foot}(g) + \dots$$

↳ For Complex Systems, prefer maximal coords

Since less hand engineering coords

- Trajectory Optimization w/ mode sequence

- Works w/ both Min/Max Coords

- Very Practical & works well

- To w/o mode sequence \Rightarrow Soft Contact

- "your mileage may vary"

\Rightarrow w/ Maximal Coords

$$\dot{x} = f(x, u, 1) \quad \theta_i(g) \geq 0$$

$\theta_i(g) \geq 0$
e.g. height above ground

$\lambda_i \geq 0$
e.g. normal force

hard contact

$$\lambda = \begin{cases} 0 & \theta_i(g) > 0 \\ -K(\theta_i(g)) & \theta_i(g) = 0 \end{cases}$$

Spring

$$\theta_i(g) \cdot \lambda_i = 0 \Rightarrow \text{e.g. either touching ground or normal force} = 0$$

$$\theta = 0$$

$$\lambda = 0$$

-Sums of Squares for Lyapunov analysis through Contact

$$M(\dot{q})\ddot{q} + C(\dot{q}, \dot{q})\dot{q} = \tilde{T}_g(\dot{q}) + Bu + \sum_i J_i(\dot{q})2; \quad \begin{matrix} \checkmark & \text{Contact Jacobian} \\ & \text{still polynomial} \end{matrix}$$

$$S = \sin \theta \quad C = \cos \theta$$

$$S^2 + C^2 = 1$$

or

Taylor expansion to make linear

$$\forall x, \lambda \quad V(g, \dot{q}) \geq 0 \quad \dot{V}(g, \dot{q}) \leq 0$$

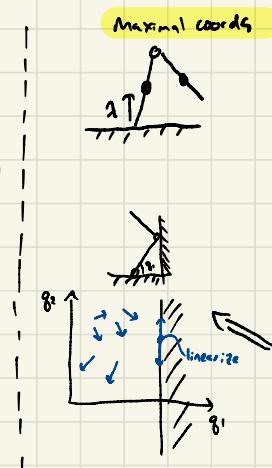
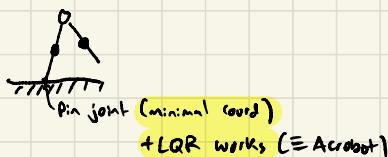
$$V(g, \dot{q}) \leq \rho \quad \lambda \geq 0 \quad \partial g \geq 0 \leftarrow S\text{-procedure / Lagrange multipliers}\$$

Chot

• SOS verification works well for simple polynomial controllers but not for e.g. MPC

-Stabilize a fixed point through contact

• Ex: Compass Gait Balancing



$$\begin{aligned} \dot{x} &= f(x, u, \lambda) \leftarrow \\ g(x, u, \lambda) &= 0 = \text{foot}(g) \quad \text{Can solve for } \lambda \text{ & insert} \\ \hookrightarrow \dot{x} &= f(x, u) \quad x \in \mathbb{R}^4 \\ &\approx Ax + Bu \end{aligned}$$

LQR(A, B, Q, R) will fail b/c not controllable
↳ intuit: linearizing around local point, loose contact behavior
↳ instead LP MPC will succeed from valid x (with $\text{foot}(g) = 0$)

Equality
Constrained LQR

$$\dot{x} = Ax + Bu \quad \text{Know additionally that } Gx = 0$$

↳ project dynamics into null space of the (linearized) constraint (smaller state space) then project back up

-Prev Solns: Stabilize for single state of contact but can't do e.g.

'push recovery' and constraint (contact) changes

↳ We are roughly stuck trying to solve across contacts

↳ Boston Dynamics does this with hand designed controllers

↳ intuit: LDL can't figure out bounce on wall to get back to fixed pt



↳ Soln: mixed integer MPC (no global sol)

- Local stabilization across modes as mixed integer optimization

$$\begin{array}{ll}\text{Min}_{x,u} & \sum Qx^T x + u^T R u \\ \text{s.t.} & x[n+1] = Ax + Bu \\ & x[0] = x_0 \\ & x[N] = x_{\text{goal}}\end{array}$$

Now:

$$x[n+1] = A^i x[n] + B^i u[n] + c^i \quad \text{when } D^i x + F^i u \leq 0$$

↳ Piecewise affine

$$\begin{array}{ll}\text{Min}_{x,u,\lambda} & \sum Qx^T x \dots\end{array}$$

$$x[n+1] = Ax[n] + Bu[n] + B\lambda[n] + C \dots$$

$$\lambda[n] \geq 0 \quad Gx[n] \geq 0$$

don't do $\lambda \cdot G \geq 0$ like before b/c not convex now

instead:

$$\begin{aligned}\lambda[n] &\leq b_i M \leftarrow \text{"big positive constant"} \\ G_i x[n] &\leq (1-b_i)M \leftarrow \text{"big M"} \\ b_i &\in \{0, 1\}\end{aligned}$$

↳ Mixed integer optimization

• Note Mixed integer always Non-Convex

• Worst case Complexity is very bad

• "Mixed-integer convex" iff ignoring integer results in convex

↳ get lower bounds \Rightarrow effective branch & bound search

↳ Very efficient solvers:

- return global optimal or infeasible

• Use case example: Atlas robot given set of foot placements determines

which foot to place in which placement at which time

- interesting: Atlas pushed against wall & using push recovery

is unreliable with full non-linear system but consistent w/ MI convex

but still only works near upright (blk linearized)

• Not realistic to run in real time

- Some methods for faster Computer:

◦ tighter formulation (from disjunctive programming)

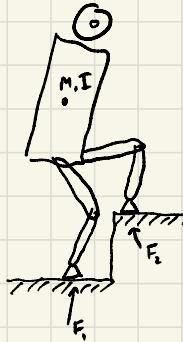
◦ Approx explicit MPC: assume contact x will be activated & find best set of initial conditions for recovering using that

◦ Lyapunov based synthesis

Lecture 16 Humanoids "Highly-articulated legged robots"

- So far: traj opt, LQR, SOS still useful w/ legs

- Today: as complexity increases so does computation, but at some point it becomes easy again



$$F_{i,z} \geq 0 \quad |F_{i,x}| \leq \mu F_{i,z}$$

- doesn't need to be massless legs

↳ can regulate COM position & total angular momentum of COM

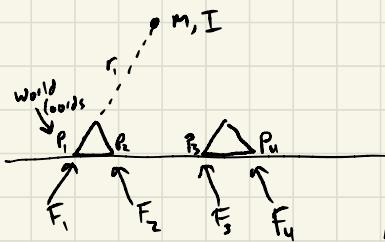
- Ignoring impact is possible:

$$\ddot{z}_{\text{COM}} = \text{const} \Rightarrow \ddot{z} = 0 \quad \& \quad \dot{\theta} = 0$$

↳ impacts impart no energy (foot speed at impact = 0)

- Flat Terrain

$$\text{LMB: } m\ddot{x} = \sum F_{i,x} \quad m\ddot{z} = \sum F_{i,z}$$



$$F_{i,z} \geq 0 \Rightarrow \ddot{z} \geq -g$$

$$|F_{i,x}| \leq \mu F_{i,z} \Rightarrow |\dot{x}| \leq \mu (\ddot{z} + g)$$

$$\text{AMB: } I\ddot{\theta} = \sum (r_{i,z} F_{i,x} - r_{i,x} F_{i,z})$$

$$r_i = p_i - \begin{bmatrix} x \\ z \end{bmatrix}$$

$$x_{\text{COM}} = \frac{\sum x_i m_i}{\sum m_i} \in \mathbb{R}^3$$

$$x_{\text{COM}} = \frac{\sum p_i x_i m_i}{\sum p_i m_i} \Rightarrow x_{\text{COM}} \in \text{Convex hull } P$$

"Support polygon"

$$\Rightarrow (m\ddot{z} + mg)(x_{\text{COM}} - x) = (z_{\text{COM}} - z)m\ddot{x} - I\ddot{\theta}$$

$$\Rightarrow \text{Assume } \ddot{\theta} = 0, \ddot{z} = 0 \Rightarrow z - z_{\text{COM}} = \text{height}$$

$$\dot{z} = 0$$

$$\Rightarrow \dot{x} = \frac{d}{h} (x - x_{\text{COM}}) \quad \text{AKA COM motion dictates what COP needs to be}$$

↳ important! in argument of force vs. position control
this says you can't control interaction forces even w/ position control (how Asimo works)

↳ works for 3D too

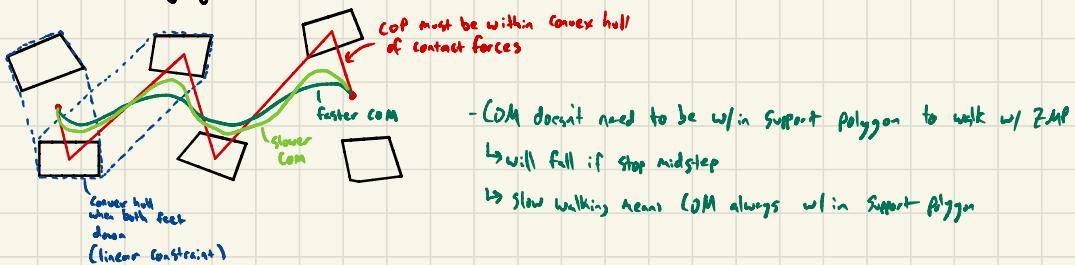
- Zero Moment Point (ZMP) : think center of pressure



- In flat ground walking X_{COP} is a zero-moment point
- requires 'light' footsteps (zero speed of foot at contact for ZMP to hold)

Three Stage Walking Plan / Control

- 1) Footstep Planning : MI Convex Optimization
- 2) Plan COM/COP : LP/QP
- 3) Fill in details (\dot{q}, \ddot{q})

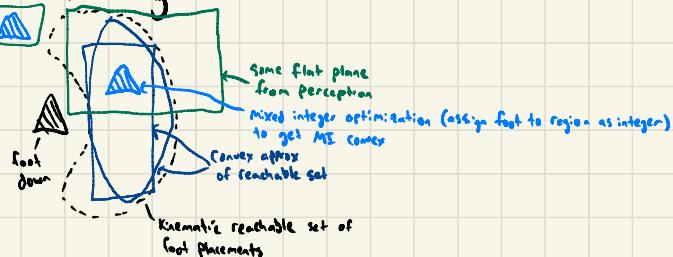


- Can use linear optimal control to jointly solve COP/COM trajectories since linear eqns & linear constraints

- Eg. give desired traj but then minimize effort via quadratic cost

- Can even use LQR to do this

Footstep Planning



- MI - Convex opt based
- Nonlinear Traj opt based
- Sampling based

} all use simplified kinematics & almost no dynamics

- Hard to do footstep planning while optimizing COM/COP because it becomes bilinear (no longer convex MI) : $\Sigma = p \times F = p_x f_x - p_z f_z \leftarrow \text{bilinear}$

↳ current research looking into this

LECTURE 17: "COMPLETE" PLANNING

- Trajectory optimization not always best for planning if NLP

- Not robust
- No guaranteed soln
- May not be optimal

- "Complete" (AI/Graph Search Community)

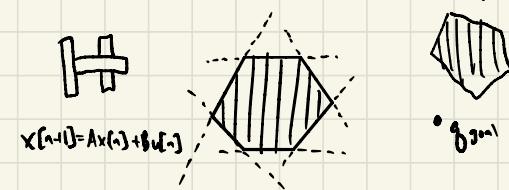
• Def: guaranteed to find a soln if one exists

◦ "Global optimal" if it finds the optimal plan

1) Decompose non-convexity in problem & search all decompositions

2) Randomized Motion Planning: very good but not great for dynamical systems

} will find something
in between



- Randomized Motion Planning "Sample-based Planning"

◦ RRT: works well for high dimensions, explores rapidly

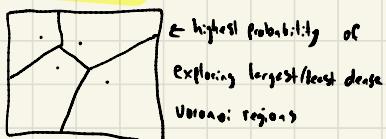


1) Random Sample

2) find closest point on current tree

3) grow tree towards sample point

◦ has a Voronoi bias



◦ Best for kinematic problems (any direction is feasible) (ignore dynamics)

◦ As $t \rightarrow \infty$: Uniform coverage of free space

◦ Even pendulum dynamic constraints make it hard to find soln

$$\min \sum (x[n] - x_{goal})^T Q (x[n] - x_{goal})$$

Interior polygon described by linear constraints

$$P \{ x \mid Ax \leq b \}$$

$$\text{Exterior: } a_1x \geq b \text{ OR } a_2x \geq b \text{ OR } a_3x \geq b$$

"breaks everything" requires Mixed Integer

↳ Disjunctive constraints w/ integer variables (binary)

$$c[n] \in \{0, 1\}$$

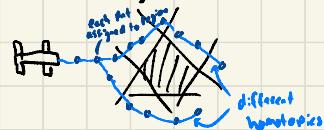
$$a_i x[n] \geq b - c[n] M$$

M: large constant

$$\sum c[n] \geq 1 \leftarrow \text{MI at least one free space occupied}$$

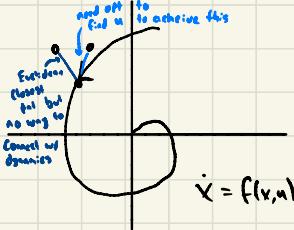
↳ MI opt uses branch and bound which uses relaxed constraints e.g. $c[n] \in [0, 1]$

↳ best when using MS cover



↳ MI can solve any NP-hard problem so takes a while

↳ Make easier (non-MI) by breaking up free space into polygons
- not complete if some free space isn't covered everything



- 1) Better sampling distributions
- 2) Better distance Metric
- 3) Better extent (e.g. Traj Optimization)

Approx

"Distance" Metrics:

- Min time to go (boundary value prob)

- Cost-to-Go

- Eg. approx by linearizing around sample & find LQR Cost-to-Go (LQR RRT*)

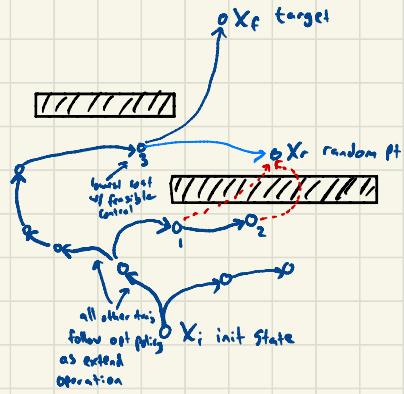
Cool! Sample based planning
for dynamical systems
"To this day one of the most
satisfying sols"

Emilio Frazzoli Method! 2002

- For real-time motion planning of helicopter

- Value iteration on grid of obstacle-free
space to get cost-to-go $J^*(\cdot, x_r)$
- ↳ use as approx of distance metric in RRT
for real-time RRT

- ↳ use pre-computed optimal policy $\Pi(\cdot, x_r)$
when no obstacle



- ↳ exploit invariants in configs to get value iteration
of 6 DoF helicopter

LECTURE 18: FEEDBACK MOTION PLANNING

- Last Time: "complete": given start & goal, will find path if it exists

- Getting to Feedback:

- Plan every timestamp (MPC)

- Roadmaps: graph across and can be reused with various start & goal

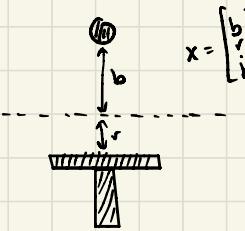
- **RRT*** (backwards trees)

- gets rid of kinks/RRT dance by constantly refining tree
- works well but can't handle disturbances

Side note: "RRT dance" funny set of intermediate configs before getting to final



- Simplest Juggling Model "Line Juggler"



$$x = \begin{bmatrix} b \\ r \\ b \end{bmatrix}$$

Assumptions

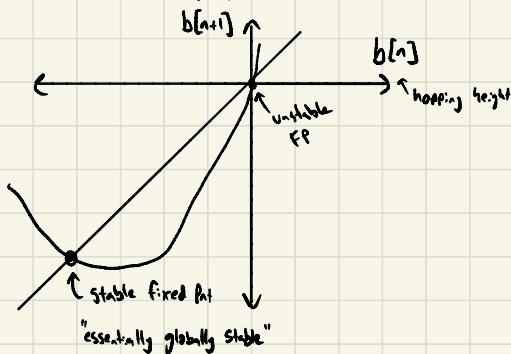
- instantaneous elastic collision
(w/ coeff of restitution)

Mirror Law:

$$\text{impose } r(t) = -kb(t)$$

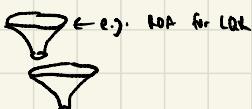
↳ collisions occur at $r=0$ $b=0$ $\dot{b}=0$

↳ one-dimensional Poincaré Map



- "Sequential Composition of Dynamically Dexterous Robot Behaviors" - R.R. Burridge

initial idea of combining Lyapunov functions from simple controllers to have R.R.A. fill state space



↳ before SOS so used experimental data to find R.R.A.s

- Probabilistic Feedback Coverage

- Lyapunov fns along trajectories

LQR Trees

$$\dot{\bar{x}}(t) = A(t)\bar{x}(t) + B(t)u \Rightarrow \bar{u}(t) = -K\bar{x}(t)$$

TV LQR

$$J(x, t) = \bar{x}^T S(t) \bar{x}$$

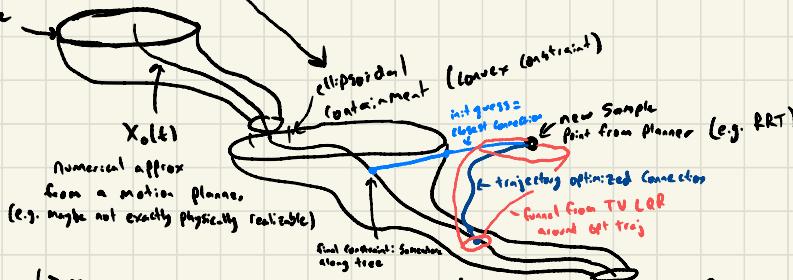
Trajectory $X_0(t)$ defined over some finite $t = [0, t_f]$

Find $P(t)$ s.t. $J(x, t) \leq P(t)$ when $J(x, t) = P(t)$ \leftarrow Note! normally say $J \leq P(t)$, instead only look at wall of funnel

^(with spine) $J(x, t_f) \leq P(t_f)$ to be inside next controller

"funnel"

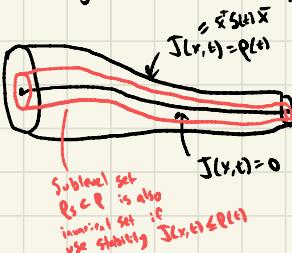
finite-time invariant set



\hookrightarrow Sparse Sampling tree because no need to sample w/in existing funnels

\hookrightarrow Think: different controllers that turn on for diff parts of state space (depending which funnel is active like a state-machine

• Issue: building RRT LQR tree is expensive \rightarrow define computation. How run online with obstacles?



$$J(x, t) \leq P(t) \rightarrow H_x \quad J(x, t) \leq P(t) \leftarrow \text{Stability}$$

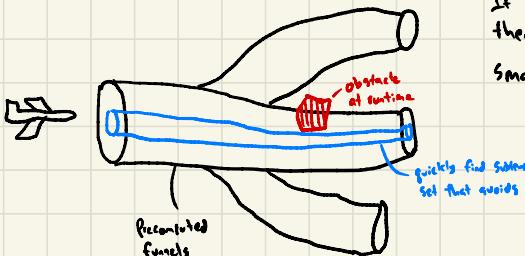
$$H_x \quad J(x, t) = P(t) \leftarrow \text{finite time invariant}$$

$$H_x \quad \alpha(t) \leq J(x, t) \leq P(t)$$

$$0 \leq \epsilon \leq 1 \leftarrow \text{avoid } 0 \& \text{ numerical tolerances}$$

If you use these stability measures instead, then can quickly shrink $P(t)$ at runtime to find smaller funnel soln. This is useful for online

Ex:



• What if you want larger ROA than given by TV LQR applied to non-linear system:

Can find new control that optimizes P

$$\dot{x} = f_1(x) + f_2(x)u \quad u = -Kx \text{ or } u = -Km(x)$$

$$\text{Given } V \quad \dot{V}(x) = \frac{\partial V}{\partial x} [f_1(x) + f_2(x)Km(x)]$$

$$\dot{V}(x) \leq 0 \text{ for all } x \quad V(x) \leq P$$

^{Monomials}

Use Sums-of-Squares Programming to search for P & K simultaneously

↳ Can start w/ K from LQR

- find V to get largest P

- look in V & find K to get even bigger P

- and alternate (Bilinear Alternations to maximize P)

↳ note: fast search for V & K that maximize P at same time



↳ LQR Trees great for Vehicles/acrobot/cartpoles/etc but NOT manipulation

because full state feedback which is not possible w/ manipulation

- shouldn't need to know full state of key-ring to pick them up

- this is reason Russ has become interested in manipulation

LECTURE 19 MODEL SYSTEMS WITH STOCHASTICITY (Robust Control)

So far $\dot{x} = f(x, u) \quad x[n+1] = f(x[n], u[n])$
 $\dot{x} = f(x, u, w) \quad x[n+1] = f(x[n], u[n], w[n])$

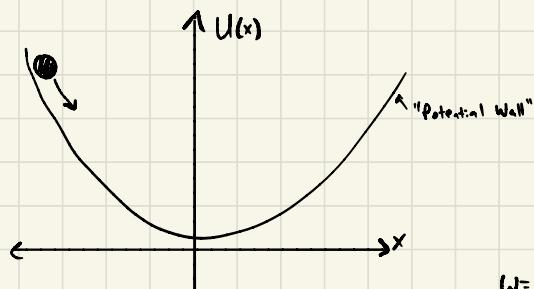
w(t) or w[n] output
of some random process

- $w[n]$

- disturbances
- Model uncertainty (process noise)
- Not measurement noise

- Additive Noise $x[n+1] = f(x[n], u[n]) + B_w w[n]$

- Ex Particle in Bowl with Brownian Motion



$$U(x) = \frac{1}{2} \alpha x^2$$

$$x[n+1] - x[n] = \left. \frac{-\partial U}{\partial x} \right|_{x[n]} + w[n]$$

$$x[n+1] = (1-\alpha)x[n] + w[n]$$

$w=0$ fixed point $x^* = 0$
 expo stable for $0 < \alpha < 2$

- $w[n]$ is Gaussian white noise

$$\begin{aligned} E[w[n]] &= 0 \\ E[w[i]w[j]] &= \begin{cases} \sigma^2 & i=j \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

- Given $x[0]$ what is $x[n]$ what is $x[\infty]$?

$x[n]$ is a random variable described by

$P_{x[n]}(\cdot)$ density function

$$P_w(w) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{w^2}{2\sigma^2}} \leftarrow \text{Gaussian}$$

$$\hookrightarrow P_{x[n+1|x[n]]}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-(1-\alpha)x)^2}{2\sigma^2}}$$

"Master Eqn"

$$P_{x[n+1]}(y) = \int_{-\infty}^{\infty} P_{x[n+1|x[n]]}(x) P_w(x) dx$$

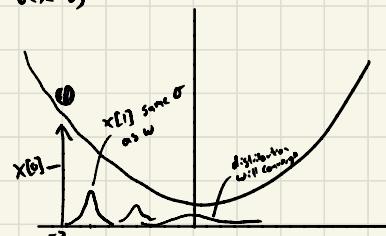
↑ convolution

output remains Gaussian

- Note init conditn $x[0]=3$

treat probability like delta fn

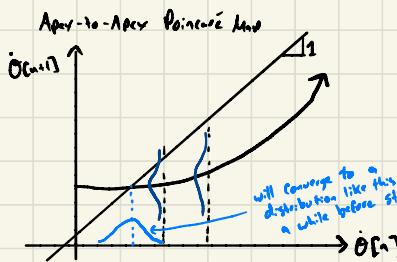
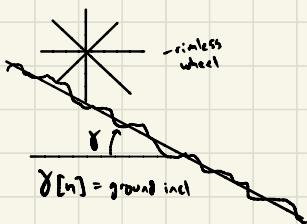
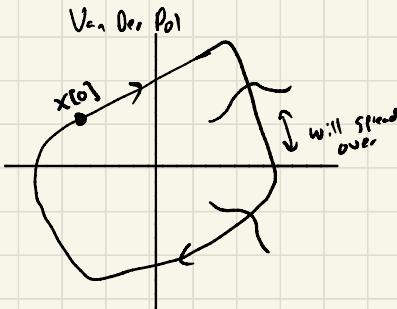
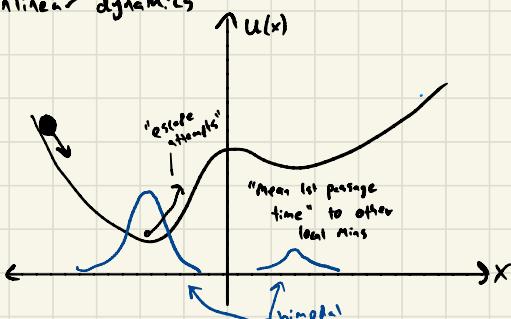
$$P_{x[0]}(x) = \delta(x-3)$$



- Steady State distribution $P_{X^*}(\cdot)$

$$\mu^* = 0 \quad (\sigma^*)^2 = \frac{\sigma^2}{2\alpha - \alpha^2}$$

- Nonlinear dynamics



- Side Note: When to worry about stochasticity?

↳ When noise distro is relatively wide compared steepness of dynamics

Stationary distribution is standing still

Since inescapable once there

↳ Think: "Any system w/ stochasticity will fill over given enough time (for long tail noise)"

↳ how to control system if final state is always fall down?

↳ "Meta-Stable" distribution - long living transients in $P_{X(n)}(\cdot)$

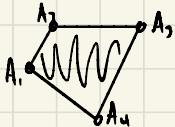
1) Worst Case design/analysis

Typically for bounded uncertainty

$P_U(\cdot)$ has finite domain (unlike Gaussian: tails $\rightarrow \infty$)

e.g. Via Common Lyapunov Function

$$\dot{x} = Ax \quad A = \sum_i \beta_i A_i \text{ affine combination}$$



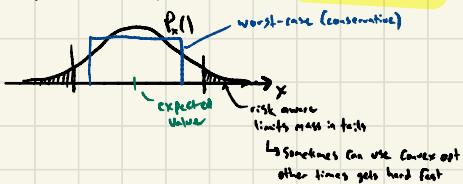
for all w , $V \leq 0$

↳ often too conservative since Considering Worst Case

2) Stochastic Control

- Write objective on $P_{x(n)}(\cdot)$ most common is expected performance

◦ Another example is Risk-Aware Control (chance constraints)



◦ Expected Value (used in Reinforcement Learning)

$$J = \sum_{n=0}^{\infty} g(x_n, u_n) \quad \text{now min } E[J] = \sum_{n=0}^{\infty} E[g(x_n, u_n)] \\ \hookrightarrow \text{now } J \text{ is random variable} \\ = \sum_{n=0}^{\infty} \left[\int_S g(x_n, u_n) p_n(x) dx \right]$$

Stochastic LQR

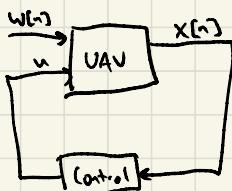
$$x[n+1] = Ax[n] + Bu[n] + B_w w[n] \\ \uparrow \text{non-zero Gaussian iid}$$

$$\min E\left[\sum x^T Q x + u^T R u\right]$$

$u = -Kx$ which is the same K in normal LQR! ↳ not same for nonlinear

$$J = \infty = x^T S x + C \\ \uparrow \text{original } S \text{ goes to } \infty \sum E[w[n]w^T[n]]$$

Ex UAV in wind

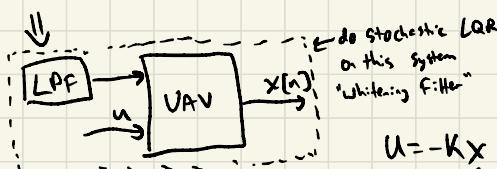


Wind is not Gaussian iid



Sidenote: LQR doesn't provide guarantees for robustness

(funny paper says there are no guaranteed stability margins for LQG)



$$u = -Kx$$

↑ UAV + wind estimation now gives good robustness

LECTURE 20 Robust Constrained Control

$$x_{t+1} = f_t(x_t, u_t) + w_t \quad \text{with constraints on } u, x$$

Disturbance $w_t \in W_t$

$$u_t = \Pi_t(x_0, \dots, x_t, u_0, \dots, u_{t-1})$$

Controller or Memory
Control Policy

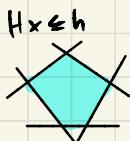
- Problem: design Control Policy Π_t such that all possible trajectories is subset of reachable trajectories Ψ
- $\Psi(x; \Pi)$
along with each run
b/c w_t

- From now on focus on linear systems

$$x_{t+1} = A_t x_t + B_t u_t + w_t \quad w_t \in W_t$$

Polytope aside:

- H-Polyhedron: define edges



$Hx \leq h$ good in high dimensions ($2n$) faces

- easy to check if within (check inequality)

| Polytope = bounded

| Polyhedron = unbounded

- U-Polytope: convex hull of vertices $\{v_1, \dots, v_n\}$

$$x = \sum_{i=1}^n \lambda_i v_i \quad \lambda_i \geq 0 \quad \sum_{i=1}^n \lambda_i = 1$$

slack



- bad in high dimensions 2^n vertices

- hard to check if in within (solve linear program)

- easy to get Minkowski sum

Minkowski Sum

$$S_1, S_2 \subset \mathbb{R}^n$$

$$S_1 \oplus S_2 = S_1 + S_2 \quad s_1 \in S_1, s_2 \in S_2$$



$X_{t+1} = A_t X_t + B_t U_t + W_t$ $W_t \in \mathbb{W}_t$ \leftarrow non-deterministic
 C_t cont. \leftarrow generator
 $P_t = \bar{X}_t + G_t Q$
 t states $t+1$ \leftarrow typically 'cube': $[-1, 1]$

Zonotope $Z(\bar{X}, G)$

↑ symmetric about center e.g.



Controller:

$$x = \bar{X}_t + G_t g \quad g \in Q$$

$$U_t = \bar{U}_t + \Theta_t g$$

↑ design parameter feedback gain

$$\begin{aligned} \text{now } A_t X_t + B_t U_t &\Rightarrow A_t (\bar{X}_t + G_t g) + B_t (\bar{U}_t + \Theta_t g) \\ &= A_t \bar{X}_t + B_t \bar{U}_t + C_t + (A_t G_t + B_t \Theta_t) g \end{aligned}$$

$$P_{t+1} = \underbrace{(A_t X_t + B_t U_t + C_t)}_{\bar{X}_{t+1}} + \underbrace{(A_t G_t + B_t \Theta_t) Q}_{G_{t+1}}$$

Linear

- Can solve w/ Trajectory Optimization (Linear Program)

$$\underbrace{\bar{X}_t \bar{U}_t G_t \Theta_t}_{\substack{\text{decision vars} \\ t=0 \dots T-1}} \quad \text{given } P_T \leq X_{goal} \quad \left. \begin{array}{l} \text{Dynamics} \\ X, u \text{ constraints: } P_t \times U_t \leq I H_0 \end{array} \right\} \begin{array}{l} \text{linear constraints} \\ \text{not shown here} \end{array}$$

T interesting note:

designing trajectory (x, u) &

feedback gain policy simultaneously (G, Θ)

- Redo \uparrow but for non-deterministic

$$x_{t+1} = A_t x_t + B_t u_t + w_t \quad w_t \in W_t = Z(\bar{w}_t, \underline{w}_t)$$

t_{start}

$$P_t = Z(\bar{x}_t, G_t)$$

$$x_t = \bar{x}_t + G_t g$$

$$u_t = \bar{u}_t + \Theta_t g$$

$$x_{t+1} = (A_t \bar{x}_t + B_t \bar{u}_t + \bar{w}_t) + \underbrace{(A_t G_t + B_t \Theta_t) g}_{\text{Minkowski sum}}$$

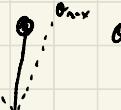
$$\begin{aligned} & \text{for boundary} \\ & Z(\bar{x}_t, G_t) \oplus Z(\bar{x}_t, G_t) \\ & = Z(\bar{x}_t + \bar{x}_t, (G_t, G_t)) \end{aligned}$$

- W_t will expand Polytopes, need to find Θ to shrink down

- same opt problem as before but Generator is longer

Robust Control

Invariance: start w/in set & need to stay within for $t \rightarrow \infty$ w/ disturbances

e.g.  $\Omega_tilde \in [\Omega_{min}, \Omega_{max}]$

- Can create LQR Tree equivalent (for linear systems & handles disturbance & constraints)

by filling up state space with Polytopes

LECTURE 21 Manipulation

(MIT 6.881: Russ' new class on manipulation)

- Guest Lecture

Dynamic Programming

- Low DoF
- $f(t), x, u, \ell(t)$ known

Stability analysis with SOS

- $f(t), x, u$ known
- $f(t)$ is vector-valued polynomial
- used for stability analysis of commercial planes

Trajectory Opt

- $\ell(t), f(t), x, u$ known
- solutions only locally opt
- used for SpaceX rocket landing

↳ falls apart for manipulation (e.g. what is state for making Sust?)

- unknown $f(t), \ell(t), g(t), x$
 $\underbrace{g(t)}_{\text{objective}}$ $\underbrace{x}_{\text{observation fn}}$

- Imitation Learning: learn policy Π from expert by observing actions & trajectories

LECTURE 22: Reinforcement Learning

- Reinforcement Learning: Collection of algorithms for "black box optimization" of stochastic optimal control problems
• Bertsekas: "ad-hoc but strangely successful" ↪ lol

black box optimization:

Only have access to $E[g(x,u)]$ costs
don't have $f(x,u)$ (model free)
how to optimize?

↳ Russ is skeptical. Should use model if have access

↳ e.g. Copter manulation hard to model

Stochastic optimal Control:

$$g(x,u) \Rightarrow E\left[\sum_{n=0}^{\infty} g(x_n, u_n)\right]$$

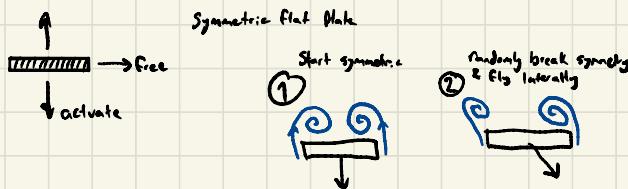
what we optimize

↳ Russ is excited about this part

↳ more generalizable

Example Fluid dynamics of flapping

The heaving foil (Jin Zhang NYU)



↳ used RL to optimize flight efficiency w/out model

↳ dynamics governed by
POE (Navier Stokes)
↳ can write $f(x,u)$
↳ CFD very expensive
↳ good use for RL

The Policy Gradient Trick (REINFORCE)

$$\min_a E[g(x)] \text{ with } x \sim p_\pi(x) \text{ Policy}$$

x is random variable that depends on parameters of trajectory

$$\begin{aligned}\frac{\partial}{\partial \alpha} E[g(x)] &= \frac{\partial}{\partial \alpha} \int dx g(x) p_\pi(x) \\ &= \int dx g(x) \frac{\partial}{\partial \alpha} p_\pi(x) \\ &= \int dx g(x) p_\pi(x) \frac{\partial}{\partial \alpha} \log p_\pi(x) \leftarrow \text{aside } y = \log u \\ &= E\left[g(x) \frac{\partial}{\partial \alpha} \log p_\pi(x)\right]\end{aligned}$$

approx w/ Monte-Carlo

$$\frac{\partial}{\partial \alpha} E\left[\sum_i^n g(x_i)\right] \approx \frac{1}{N} \sum_{n=0}^N g(x_n) \frac{\partial}{\partial \alpha} \log p_\pi(x_n)$$

↑ gives direction to update parameters or
after running many trials

$$\begin{aligned}y &= \log u \\ \frac{\partial y}{\partial \alpha} &= \frac{1}{u} \frac{\partial u}{\partial \alpha}\end{aligned}$$

$$\frac{\partial}{\partial \alpha} \log p_\pi(x) = \frac{1}{p_\pi(x)} \frac{\partial}{\partial \alpha} p_\pi(x)$$

Caveat: true of expected value is
a weak statement

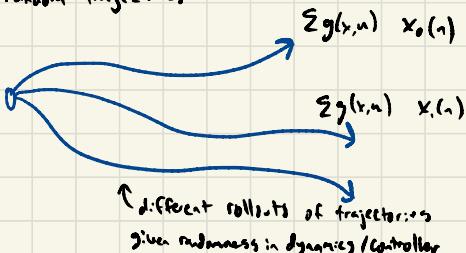
↳ Think: Surprising! gradient of performance depends on cost & gradient of controller but not on dynamics/initial conditions

Rewritten in slides as:

$$\frac{\partial}{\partial \alpha} E \left[\sum_{t=0}^T g(x[t], u[t]) \right] = E \left[\sum_{t=0}^T g(x[t], u[t]) \underbrace{\frac{\partial}{\partial \alpha} \log p_\pi(u[t] | x[t])}_{\text{e.g. last on stage #3 only influenced by stages 0:3}} \right]$$

Intuition:

N random trajectories



e.g. Simple controller $u = -Kx + w$

↑ random number in Policy helps explore space

↑ update based on Policy gradient

update K so that rollouts with w giving low costs

are more likely in the future

↳ inefficient update: many Monte-Carlo rollouts to get good policy

↳ Same update can be used for simple spline parameterization of trajectory or deep neural network

Black-Box Optimization:

$$\min_{\alpha} g(\alpha)$$

How to do gradient-free optimization?

- Idea: finite differences

$$\frac{\partial g}{\partial \alpha} \Big|_{\alpha} \approx \left[\frac{g(\alpha + \epsilon) - g(\alpha)}{\epsilon} \right]$$

↑
n parameters ↑
n+1 evaluations ↪ Can be expensive

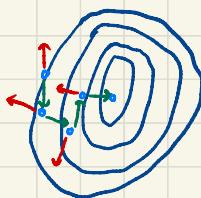
$$\text{Then do gradient descent: } \alpha[k+1] = \alpha[k] - \eta \frac{\partial g}{\partial \alpha} \Big|_{\alpha[k]}$$

Stochastic gradient descent:

- as long as going "downhill" can get away w/ less eval of $g(\alpha)$

- Simpler idea: "Weight Perturbation"

$$\Delta \alpha = -\eta [g(\alpha + \beta) - g(\alpha)] \beta$$



↑ performs stochastic gradient descent

↑ choose 1 of 2 directions

$$E[\Delta\alpha] = -\gamma \sigma^2 \frac{\partial g}{\partial \alpha}$$

↳ expected avg is in direction of actual gradient

- Often wins over finite difference over it's not to reach local min

- similar idea: "Trust Region Optimization" bound step size

- Even simpler estimate

$$\Delta\alpha = -\frac{\gamma}{\sigma^2} [g(\alpha + \delta) - b] \beta$$

↳ baseline "expected performance"

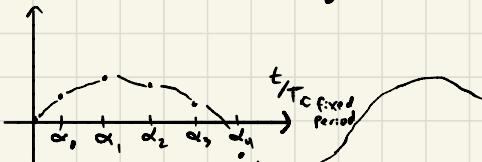
- still stochastic gradient descent: $E[\Delta\alpha] \approx -\frac{\partial g}{\partial \alpha}$

- cool: simple algo w/ strong property

- uncool: this is a weak property (same property as policy gradient trick)

Key limitation of RL today: "Sample Complexity"

Example Policy params for hearing fail



↳ better than parameterizing by Fourier coeffs: better to have all params contribute about the same to cost

- Cost: More inverse cost of transport

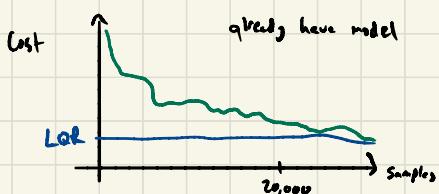
- converge consistently to triangle wave



"extremum seeking control" "iterative learning control"

- works well for repetitive experiments

- Example when not to use RL: Quadrotor stabilization



- Given accurate sim & lots of sim yrs can do robot head manipulation
- Takeaway: current state takes too many trials. Only marginally better than running billions of random trials and taking best one
- still use model when available

LECTURE 23: REINFORCEMENT LEARNING II

- RL is blackbox optimization for stochastic optimal control

- Last time: Policy Gradient Methods

$$u = \Pi_\alpha(x) \quad x = f(x, u) \quad \text{e.g.}$$

α parameters
search for opt. parameters

$$y = h(x, u)$$

$$\bullet u = -Kx$$

$$\bullet u = \text{DNN}(x)$$

$$\bullet u = -\Pi_\alpha(y) \quad \text{output} \quad \text{"output feedback"}$$

↳ powerful, unlike e.g. LQR: full-state feedback

↳ gradient descent: trial & error \rightarrow slow

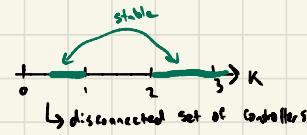
• Two Main Concerns

1) Sample complexity: e.g. 20,000 iterations for gradient

2) Non-Convexity

$$\begin{aligned} \dot{x} &= Ax + Bu & A = \begin{bmatrix} 0 & 0 & 2 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} & B = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \\ y &= Cx & (= [1 \ 1 \ 3]) \end{aligned}$$

$$u = -K_j$$



↳ Intuit: policy gradient search for K may fail to jump to find global min

↳ deep NN w/ many params can overcome local min by always finding new direction w/ lower cost

- Today - Q-Learning

- Actor-Critic

Idea: Learning Value Functions

Policy Gradient
inefficient

Dynamic Programming
efficient

↳ eval dynamics exactly
once from terminal conditions
back to $t=0$

↳ requires model

$$J^*(x)$$

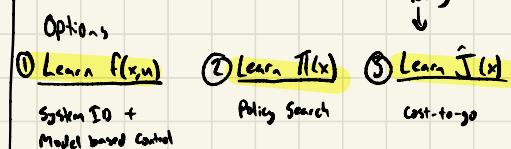
↳ optimal

$$J^*(x)$$

↳ for some policy

$$\hat{J}(x)$$

↳ approx



$$J(x) = E \left[\sum_{n=0}^{\infty} \gamma^n g(x[n], u[n]) \right]$$

↗ discount
 ↗ $0 < \gamma < 1$
 ↗ keep bounded
 ↗ value never in time more

$$x(n+1) = f(x(n), u(n), w(n))$$

- Learn only $\hat{J}^*(x)$:

$$\hat{J}^*(x) = \min_u [g(x, u) + \gamma E[\hat{J}^*(f(x, u))]]$$

$$\pi^*(x) = \arg\max_u [\dots]$$

↳ given cost-to-go also have policy (if know $f(x, u)$)

not available in RL model-free

- Two Work Arounds (not having $f(x, u)$ to map $J^* \rightarrow \pi^*$)

1) Q-Learning $Q(x, u)$ instead of $J(x)$

2) Actor-Critic: Learn $\pi(x)$ & $J^*(x)$ Simultaneously

- Q-functions

• Define $Q^*(x, u) = E[g(x, u) + J^*(f(x, u))]$ \Rightarrow think: cost of single step control + remaining cost-to-go

$$J^*(x) = \min_u Q^*(x, u)$$

J has dim(u) inputs (scalar output)

Q has $\dim(x) + \dim(u)$ inputs (scalar outputs)

◦ Standard Method: Q is a Deep Neural Network (DNN)

- QT-Opt says routine non-linear opt of $Q^*(x, u)$

- e.g. decent robot arm form

- Actor-Critic

$$\text{policy gradient } \Delta \alpha = -\gamma [J(x+\beta) - b] \underbrace{\left[\frac{\partial}{\partial \alpha} \log f_{\pi(x)}(x|\alpha) \right]^T}_{\text{baseline}} \quad \leftarrow \text{similar from last lecture: } [g(\alpha+\beta) - b]^T$$

\uparrow one time experiment

Idea: use $J(x)$ as the baseline expected value (better baseline than running an extra trial)
"critic"

- Two time scale - Controller should change fast but cost fn change slowly

- if reset controller (policy) but keep Value function (cost-to-go), controller converges quickly

◦ indicates importance of learning cost-to-go & not just control policy (Policy gradient)

- Russ thinks this method of learning Value fn is most promising

◦ E.g. AlphaGo did this (though they also used a model & full state info)

- use Monte-Carlo tree search (Stochastic MPC) where terminal cost

is cost-to-go learned by the value function

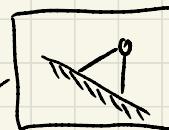
-Course Summary:



- Dynamic Programming / Value iteration
- Simple non-linear system w/ few DoF



- Traj Opt
- LQR
- SOS
- LQR-Trees

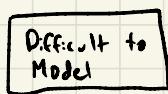


- Hybrid
- Limit Cycles

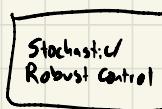


Humanoids + Manipulation

- Higher DoF
- Sample Motion planning
- MI - Convex Programming



RL



◦ Polytopic Uncertainty