

Planning using dynamic optimization

© Chris Atkeson 2004

Problem characteristics

- Want optimal plan, not just feasible plan
- We will minimize a cost function $C(\text{execution})$. Some examples:
- $C() = c_T(x_T) + \sum c(x_k, u_k)$: deterministic with explicit terminal cost function
- $C() = E(c_T(x_T) + \sum c(x_k, u_k))$: stochastic

Examples

- A number of us are currently working on humanoid locomotion. We would like the humanoid to be able to walk, run, vary speed, turn, sit, get up from a chair, handle steps, kick a ball, avoid obstacles, handle rough terrain, ... [movies]
- The next assignment will be to write a controller for a marble maze game.

Dynamic Optimization

- General methodology is dynamic programming (DP).
- We will talk about ways to apply DP.
- Requirement to represent all states, and consider all actions from each state, lead to "curse of dimensionality": $R_x^d R_u^d$
- We will talk about special purpose solution methods.

Dynamic Optimization Issues

- Discrete vs. continuous states and actions?
- Discrete vs. continuous time?
- Globally optimal?
- Stochastic vs. deterministic?
- Clocked vs. autonomous?
- What should be optimized, anyway?

Policies vs. Trajectories

- $u(t)$ open loop trajectory control
- $u = u_{ff}(t) + K(x - x_d(t))$ closed loop trajectory control
- $u(x)$ policy

Types of tasks

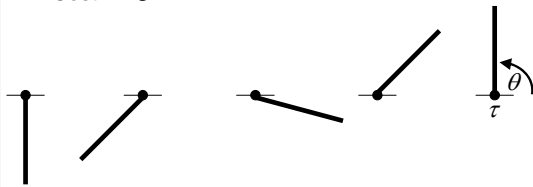
- Regulator tasks: want to stay at x_d
- Trajectory tasks: go from A to B in time T, or attain goal set G
- Periodic tasks: cyclic behavior such as walking

Typical reward functions

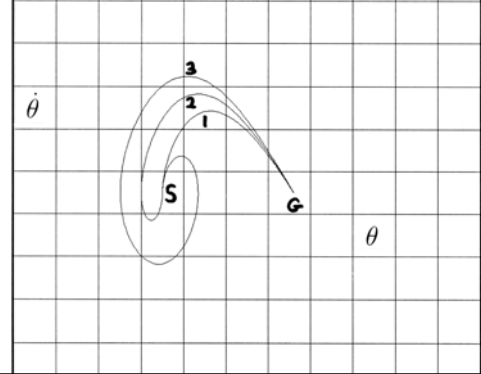
- Minimize error
- Minimum time
- Minimize tradeoff of error and effort

Example: Pendulum Swingup

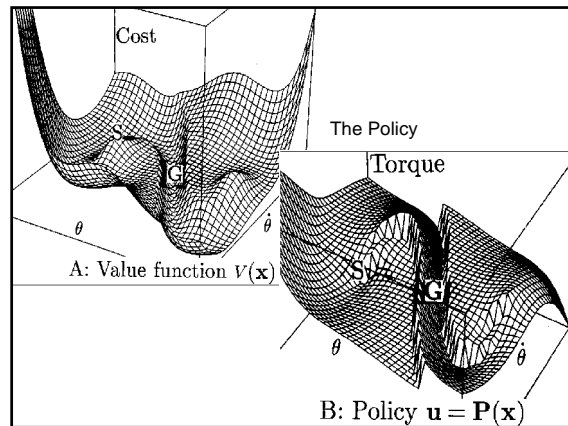
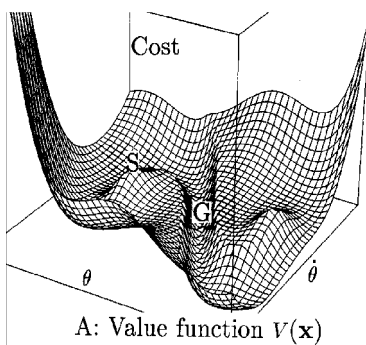
- State: $\mathbf{x} = (\theta, \dot{\theta})$
- Action: $\mathbf{u} = (\tau)$
- Cost: $\mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u}$



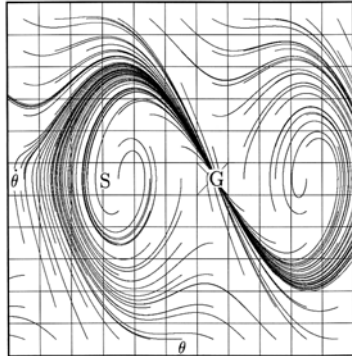
Possible Trajectories



Global Planning Using Dynamic Programming



Trajectories



Discrete Dynamic Programming

How to do discrete deterministic DP (specified time)

- Dynamics: $x_{k+1} = f(x_k, u_k)$
- Cost: $C() = c_T(x_T) + \sum c(x_k, u_k)$
- Value function $V_k(x)$ is represented by table.
- $V_T(x) = c_T(x)$
- For each x , $V_k(x) = \min_u (c(x, u) + V_{k+1}(f(x, u)))$
- This is Bellman's Equation
- This version of DP is value iteration
- Can also tabulate policy: $u = \pi_k(x)$

How to do discrete deterministic DP (no specified time)

- Cost: $C() = \sum c(x_k, u_k)$
- $V_N(x)$ is a guess, or all zeros.
- Apply Bellman's equation.
- $V(x)$ is given by $V_k(x)$ when V stops changing.
- Goal needs to have zero cost, or need to discount so $V()$ does not grow to infinity:
- $V_k(x) = \min_u (c(x, u) + \gamma V_{k+1}(f(x, u)))$, $\gamma < 1$

Discrete Policy Iteration

- $u = \pi(x)$: general policy (a table in discrete case).
- *) Compute $V^\pi(x)$:

$$V_k^\pi(x) = c(x, \pi(x)) + V_{k+1}^\pi(f(x, \pi(x)))$$
- Update policy $\pi(x) = \operatorname{argmin}_u (c(x, u) + V^\pi(f(x, u)))$
- Goto *)

Discrete Stochastic DP

- Cost: $C() = \sum E(c(x_k, u_k))$
- Bellman's equation now involves expectations:
- $V_k(x) = \min_u E(c(x, u) + V_{k+1}(f(x, u)))$

$$= \min_u (c(x, u) + \sum p(x_{k+1}) V_{k+1}(x_{k+1}))$$
- Modified Bellman's equation applies to value and policy iteration.
- May need to add discount factor.

Discrete DP will work for Maze Assignment

- Can have integral states and actions, and measure time in steps, so:
- $\text{pos}_{k+1} = \text{pos}_k + \text{vel}_k$
- $\text{vel}_{k+1} = \text{vel}_k + \text{acc}_k$
- Ball has linear dynamics, except at collisions
- Discrete DP has problems with nonlinear dynamics

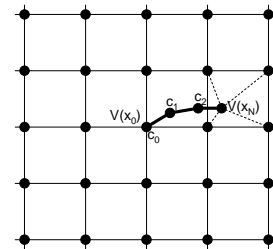
Continuous State/Action DP

- Time is still discrete.

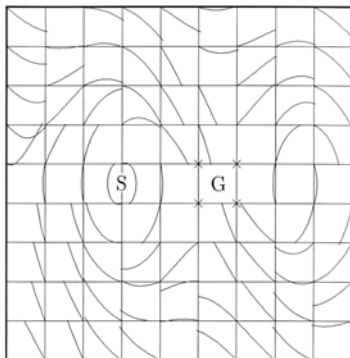
How to handle continuous states and actions (value iteration)

- Discretize value function $V()$
- At each V point (x_0), generate trajectory segment of length N by minimizing $C(u) = \sum c(x_k, u_k) + V(x_N)$
- $V(x_N)$: interpolate surrounding $V()$
- N typically determined by when $V(x_N)$ independent of $V(x_0)$
- Use favorite continuous function optimizer to search for best u when minimizing $C(u)$
- Update $V()$ at that cell.

State Increment Dynamic Programming (Larson)

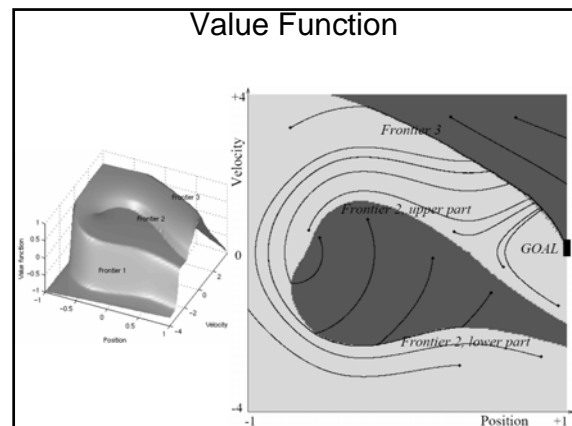
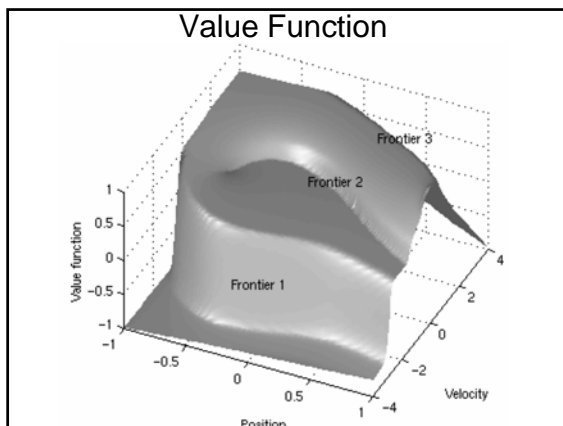
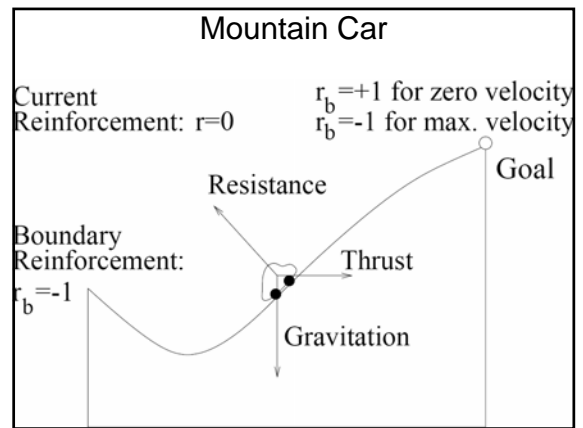
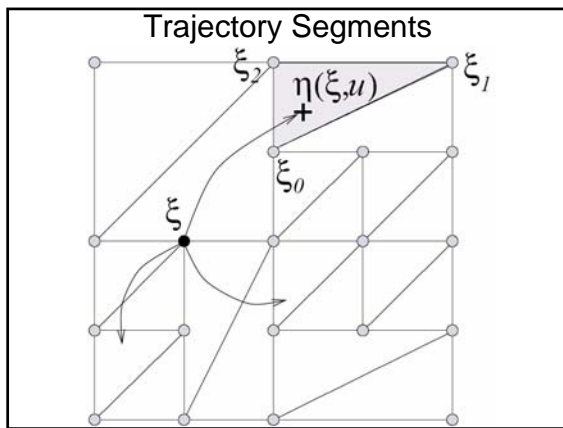
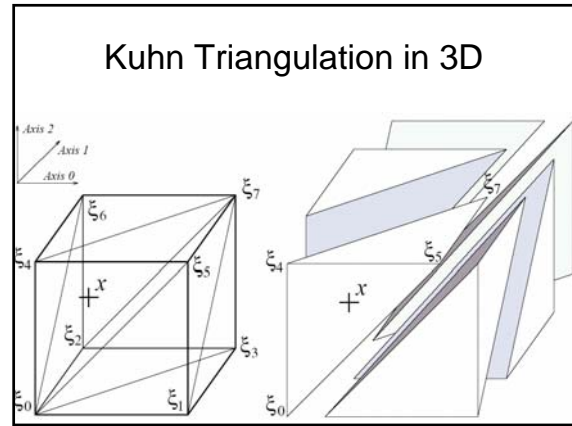
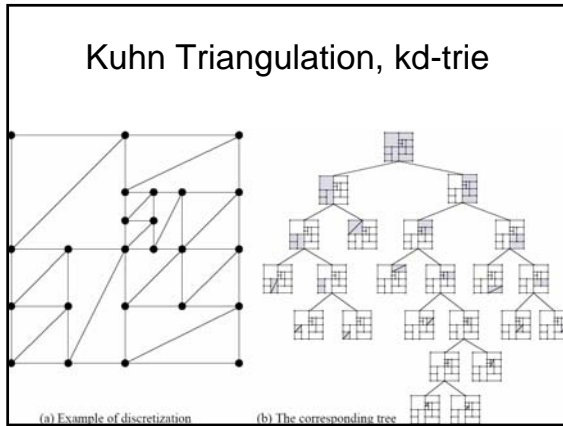


State Increment Dynamic Programming

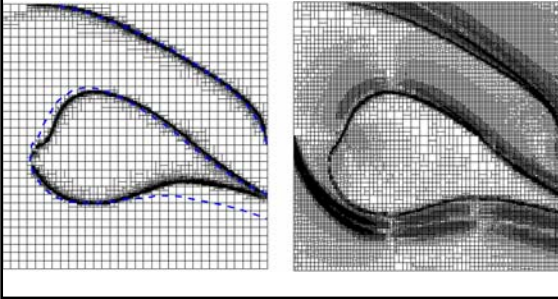


Munos and Moore, Variable Resolution Discretization in Optimal Control

Machine Learning, 49 (2/3),
291-323, 2002



Discretizations



Policy Iteration: Continuous x, u

- Discretize policy:
- Each cell in table has constant u , or
- u as knot points for linear or higher order spline
- *) Same kind of trajectory segments used to compute $V_k^\pi(x) = \sum c(x, \pi(x)) + V_{k+1}^\pi(x_N)$
- Optimize policy $\pi(x) = \operatorname{argmin}_u (c(x, u) + V^\pi(f(x, u)))$ using favorite continuous function optimizer.
- Goto *)

Stochastic DP: Continuous x, u

- Cost: $C() = \sum E(c(x_k, u_k))$
- Do Monte Carlo sampling of process noise for each trajectory segment (many trajectory segments), or
- Propagate analytic distribution (see Kalman filter)
- Bellman's equation involves expectations:
- $V_k(x) = \min_u E(c(x, u) + V_{k+1}(f(x, u)))$

Regulator tasks

- Examples: balance a pole, move at a constant velocity
- A reasonable starting point is a Linear Quadratic Regulator (LQR controller)
- Might have nonlinear dynamics $x_{k+1} = f(x_k, u_k)$, but since stay around x_d , can locally linearize $x_{k+1} = Ax_k + Bu_k$
- Might have complex scoring function $c(x, u)$, but can locally approximate with a quadratic model $c \approx x^T Qx + u^T Ru$
- `dlqr()` in matlab

LQR Derivation

- Assume $V()$ quadratic: $V_{k+1}(x) = x^T V_{xx;k+1} x$
- $C(x, u) = x^T Qx + u^T Ru + (Ax + Bu)^T V_{xx;k+1} (Ax + Bu)$
- Want $\partial C / \partial u = 0$
- $B^T V_{xx;k+1} Ax = (B^T V_{xx;k+1} B + R)u$
- $u = Kx$ (linear controller)
- $K = -(B^T V_{xx;k+1} B + R)^{-1} B^T V_{xx;k+1} A$
- $V_{xx;k} = A^T V_{xx;k+1} A + Q + A^T V_{xx;k+1} BK$

More general LQR equations

$$x_{k+1} = f(x, u) \approx Ax + Bu + c$$

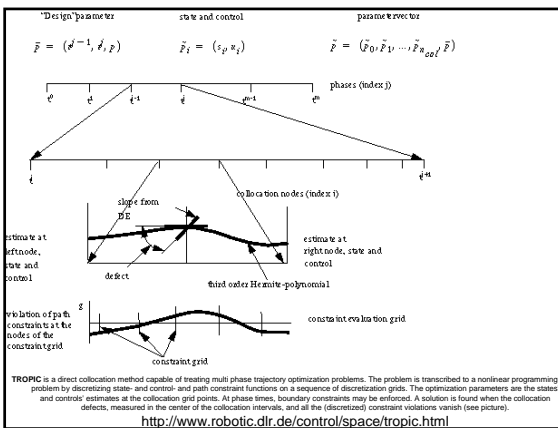
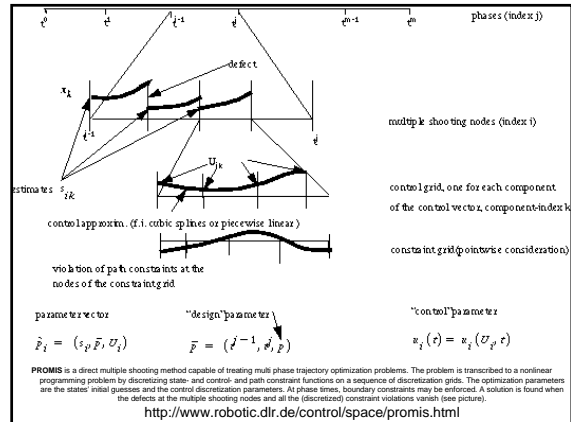
$$L(x, u) \approx \frac{1}{2} x^T Q x + \frac{1}{2} u^T R u + x^T S u + t^T u$$

$$V(x) \approx V_0 + V_x x + \frac{1}{2} x^T V_{xx} x$$

$$u^{opt} = -(R + B^T V_{xx} B)^{-1} x (B^T V_{xx} A x + S^T x + B^T V_{xx} c + V_x B + t)$$

Trajectory Optimization (open loop)

- Calculus of variations
- Multiple shooting
- Function optimization
 - Represent $x(t)$ and $u(t)$ as splines, knot point vector θ
 - Optimize $\text{cost}(\theta)$ with dynamics $x_{k+1}=f(x_k, u_k)$ a constraint or with dynamic error part of cost.
 - DIRCOL example of current state of the art.



Trajectory Optimization (closed loop)

- Differential Dynamic Programming (local approach to DP).

Propagate Value Function $V()$ Along Trajectories

$$V(x) \approx V_0 + V_x x + \frac{1}{2} x^T V_{xx} x$$

$$Z_x = V_x A + Q(x - x_d)$$

$$Z_u = V_x B + R(u - u_d)$$

$$Z_{xx} = A^T V_{xx} A + Q$$

$$Z_{ux} = B^T V_{xx} A + S$$

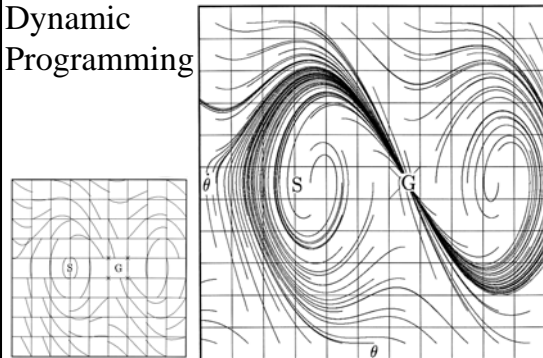
$$Z_{uu} = B^T V_{xx} B + R$$

$$K = Z_{uu}^{-1} Z_{ux}$$

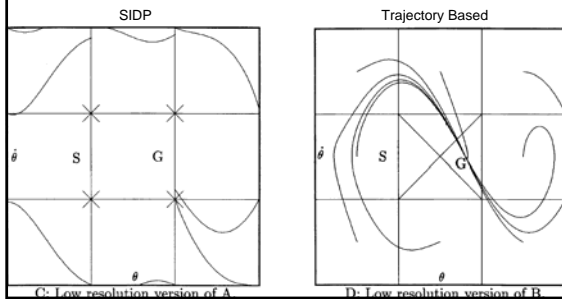
$$V_{x_{k-1}} = Z_x - Z_u K$$

$$V_{xx_{k-1}} = Z_{xx} - Z_{xu} K$$

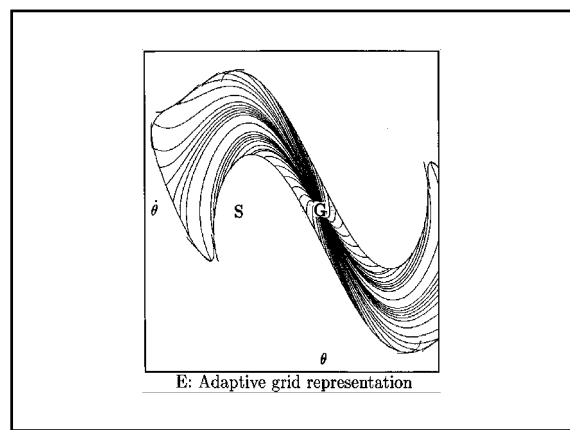
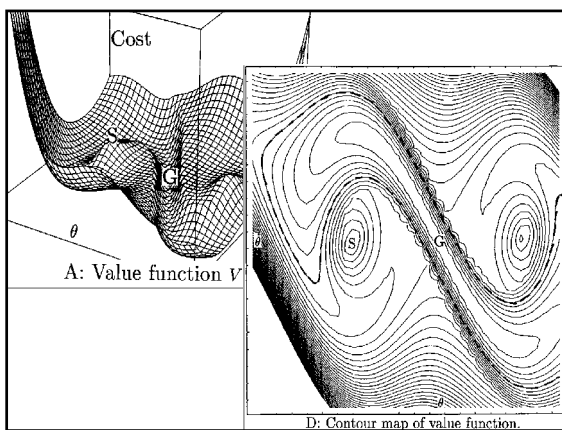
Trajectory-Based Dynamic Programming



Full Trajectories Helps Reduce Resolution Needed



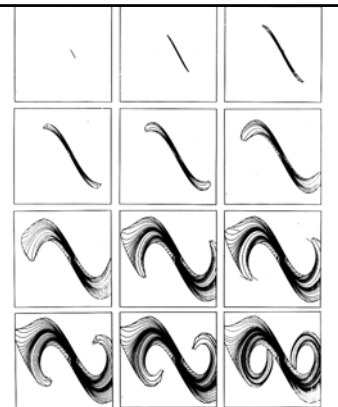
An Adaptive Grid Approach



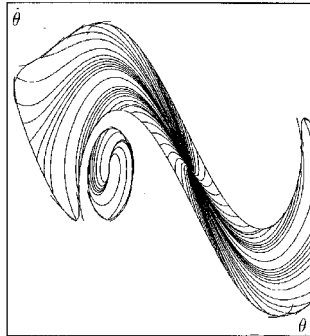
Global Planning Propagate Value Function Across Trajectories in Adaptive Grid

$$V_{01} \approx V_{02} + V_{x2}(\mathbf{x}_1 - \mathbf{x}_2) + \frac{1}{2}(\mathbf{x}_1 - \mathbf{x}_2)^T V_{xx2}(\mathbf{x}_1 - \mathbf{x}_2)$$

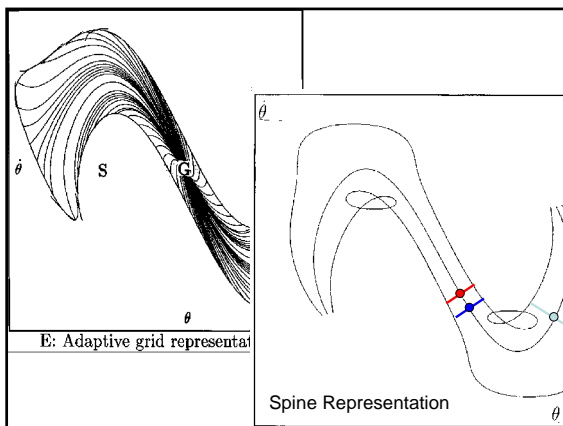
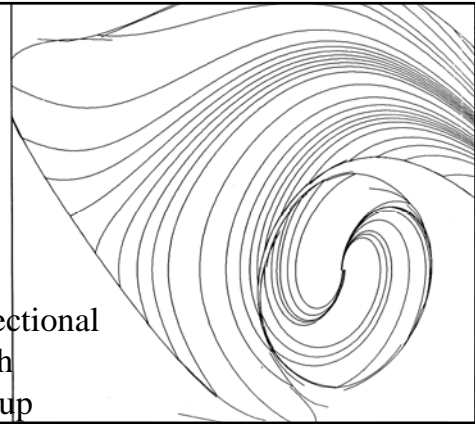
Growing the Explored Region: Adaptive Grids



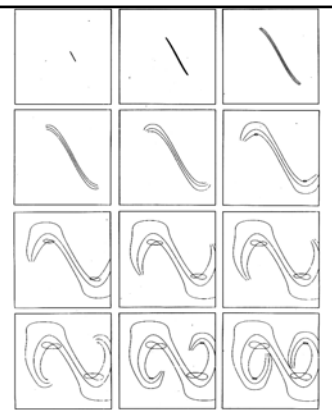
Bidirectional Search



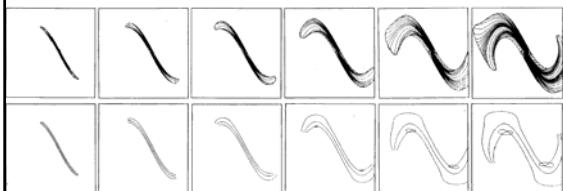
Bidirectional Search Closeup



Growing the Explored Region: Spine Representation

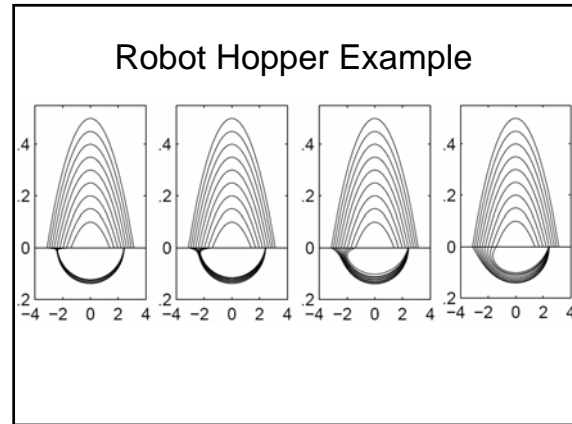
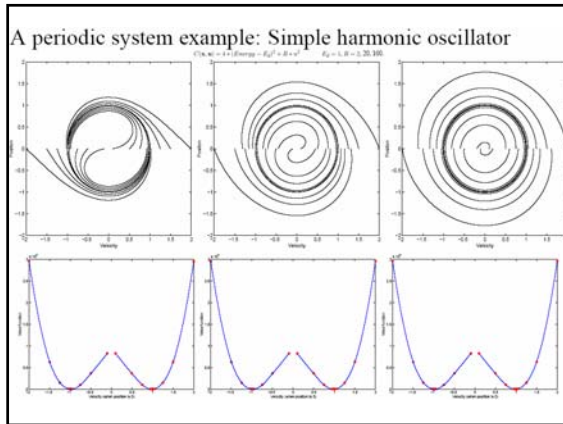


Comparison



What Changes When Task Periodic?

- Discount factor means $V()$ might increase along trajectory. $V()$ cannot always decrease in periodic tasks.



Policy Search

- Parameterized policy $u = \pi(x, \theta)$, θ is vector of adjustable parameters.
- Simplest approach: Run it for a while, and measure total cost.
- Use favorite function optimization approach to search for best θ .
- There are tricks to improve policy comparison, such as using the same perturbations in different trials, and terminating trial early if really bad (racing algorithms).

Policy Search For Structured Policies: Gradient Descent

$$J(\theta) = \int_{\mathbf{x}_0} p(\mathbf{x}_0) V^{\pi}(\mathbf{x}_0, \theta) d\mathbf{x}_0 \approx \sum_{\mathbf{x}_0} p(\mathbf{x}_0) V^{\pi}(\mathbf{x}_0, \theta)$$

$$\nabla J(\theta) \approx \sum_{\mathbf{x}_0} p(\mathbf{x}_0) \frac{\partial V^{\pi}(\mathbf{x}_0, \theta)}{\partial \theta}$$

Computing the derivatives of $V()$

$$V^{\pi}(\mathbf{x}_k, \theta) = r(\mathbf{x}_k, \pi(\mathbf{x}_k, \theta)) + \lambda V^{\pi}(\mathbf{x}_{k+1}, \theta)$$

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \pi(\mathbf{x}_k, \theta))$$

$$V = r + \lambda V^{k+1}$$

$$V_{\theta} = r_{\mathbf{u}} \pi_{\theta} + \lambda (V_{\mathbf{x}}^{k+1} \mathbf{f}_{\mathbf{u}} \pi_{\theta} + V_{\theta}^{k+1})$$

$$V_{\mathbf{x}} = r_{\mathbf{x}} + r_{\mathbf{u}} \pi_{\mathbf{x}} + \lambda (V_{\mathbf{x}}^{k+1} \mathbf{f}_{\mathbf{x}} + V_{\mathbf{x}}^{k+1} \mathbf{f}_{\mathbf{u}} \pi_{\mathbf{x}})$$

Policy Search: Stochastic Case

$$J(\theta) \approx \mathbb{E} \left(\sum_{\mathbf{x}_0} p(\mathbf{x}_0) V^{\pi}(\mathbf{x}_0, \theta) \right) = \sum_{\mathbf{x}_0} p(\mathbf{x}_0) \mathbb{E} (V^{\pi}(\mathbf{x}_0, \theta))$$

$$\mathbb{E} (V^{\pi}(\mathbf{x}_0, \theta)) \approx \sum_{k=0}^T \lambda^k (r(\hat{\mathbf{x}}_k, \hat{\mathbf{u}}_k) + \text{Trace}(\Sigma(k)(r_{\mathbf{xx}} + r_{\mathbf{xu}} \pi_{\mathbf{x}} + \pi_{\mathbf{x}}^T r_{\mathbf{ux}} + \pi_{\mathbf{x}}^T r_{\mathbf{uu}} \pi_{\mathbf{x}}) | \hat{\mathbf{x}}_k, \hat{\mathbf{u}}_k, \theta))$$

$$\approx J_d + \sum_{k=0}^T \lambda^k \text{Trace}(\Sigma(k) \mathbf{R}_{\mathbf{xx}}(k))$$

Partially Observable Markov Decision Processes (POMDPs)

- Plan using belief state (too expensive?)
- Certainty equivalent approaches: use maximum likelihood estimate of state.
- Policy search
- Dual control problem: want to control, but also want to perturb to reduce uncertainty.

Planning For Dynamic Tasks

- The computational cost of planning is the big challenge for model-based RL.
- Local planning is fast, but only locally optimal.
- Global planning is expensive, but globally optimal.
- Can we combine local and global planning to get fast planning with good plans?

How to do marble maze task: Solving one maze

- Path plan, then LQR servo: A*, RRT, PRM
- Potential field in configuration space.
- Potential field in state space.
- A*/DP in discretized state space.
- Continuous state/action DP
- Policy search

But what can you learn that generalizes across mazes?

Planning and Learning

- Learn better model, and replan.
- Plan faster
 - Initialize value function or policy
 - Find best meta-parameters
 - Find best planning method
- Make better plans
 - Find better optima
 - More robust plans (plan for modeling error)