

# Optimal and Learning Control for Autonomous Robots Lecture 10



Jonas Buchli/ Farbod Farshidian  
Agile & Dexterous Robotics Lab

# Reading

- Peters, Jan, and Stefan Schaal. "Reinforcement learning of motor skills with policy gradients."
- Deisenroth, Marc Peter, Gerhard Neumann, and Jan Peters. "A Survey on Policy Search for Robotics." (2013). [Section 2.2]



# Outline

- Natural Gradient
- episodic Natural Actor Critic (eNAC)

# Policy Gradient Theorem (PGT)

- Gradient in Policy Gradient Theorem (PGT)

$$\nabla_{\theta}^{PGT} J(\theta) = E_{p_{\theta}(\tau)} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} (\log \pi_{\theta}(u_t | x_t)) (Q_t^{\pi}(x_t, u_t) - b_t) \right]$$

- If  $b_t = V_t^{\pi}(x_t)$

$$\nabla_{\theta}^{PGT} J(\theta) = E_{p_{\theta}(\tau)} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} (\log \pi_{\theta}(u_t | x_t)) (Q_t^{\pi}(x_t, u_t) - V_t^{\pi}(x_t)) \right]$$

Advantage function  $A_t^{\pi}(x_t, u_t) = Q_t^{\pi}(x_t, u_t) - V_t^{\pi}(x_t)$

$$\nabla_{\theta}^{PGT} J(\theta) = E_{p_{\theta}(\tau)} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} (\log \pi_{\theta}(u_t | x_t)) A_t^{\pi}(x_t, u_t) \right]$$



# Unbiased estimation of gradient

- We need to approximate the advantage function

$$\nabla_{\theta}^{PGT} J(\theta) = E_{p_{\theta}(\tau)} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} (\log \pi_{\theta}(u_t | x_t)) A_t^{\pi}(x_t, u_t) \right]$$

- The approximation is done through function approximation
- This function approximation should not cause bias in the gradient estimation
- **But** every function approximation has error
- **So** the error should be orthogonal to the gradient direction



# Compatible function approximation

- The function approximation for the advantage function should minimize the expectation of the square error (ESE)

$$\min_w E_{p_\theta} \left[ \left( A_t^\pi(x_t, u_t) - f_w(x_t, u_t) \right)^2 \right]$$

- The function approximation is linear with respect to its parameters

$$f_w(x_t, u_t) = w^T \underbrace{\nabla_\theta (\log \pi_\theta(u_t | x_t))}_{\text{Base functions are gradient of policy}}$$

Base functions are  
gradient of policy

- $w$  should be found through minimizing the ESE



# PGT with compatible function approximation

- Using the function approximation in the PGT gradient while the baseline is taken as value function

$$\nabla_{\theta}^{PGT} J(\theta) = G_{\theta} w$$

$$G_{\theta} = E_{p_{\theta}(\tau)} \left[ \sum_{t=0}^T \nabla_{\theta} (\log \pi_{\theta}(u_t | x_t)) \nabla_{\theta} (\log \pi_{\theta}(u_t | x_t))^T \right]$$

# Goals of Natural Gradient

- Avoiding quick decrease in the exploration ability
- keeping the exploitation of the gradient information local

## Idea of Natural Gradient

limit the changes of the policy distribution or equivalently the changes of trajectory distribution





# Natural Gradient

So the problem statement is:

Find the parameter change which maximize the cost function below while keeping distance between two distributions  $\varepsilon$

$$\max_{\Delta\theta} J(\theta + \Delta\theta) \approx J(\theta) + \Delta\theta^T \nabla_{\theta} J$$

$$s.t. \quad \varepsilon = d_{KL}(p_{\theta}(\tau) \parallel p_{\theta+\Delta\theta}(\tau)) \approx \frac{1}{2} \Delta\theta^T F_{\theta} \Delta\theta$$



Learning rate

$$\Delta\theta = \alpha_n F_{\theta}^{-1} \nabla_{\theta} J$$

$$F_{\theta} = E_{p_{\theta}(\tau)} \left[ \sum_{t=0}^T \nabla_{\theta} (\log \pi_{\theta}(u_t | x_t)) \nabla_{\theta} (\log \pi_{\theta}(u_t | x_t))^T \right]$$



# Natural Policy Gradient

Using PGT gradient in natural gradient format

$$\begin{array}{l}
 \nabla_{\theta}^{NG} J(\theta) = F_{\theta}^{-1} \nabla_{\theta}^{PG} J(\theta) \\
 \nabla_{\theta}^{PGT} J(\theta) = G_{\theta} w
 \end{array}
 \left. \vphantom{\begin{array}{l} \nabla_{\theta}^{NG} J(\theta) = F_{\theta}^{-1} \nabla_{\theta}^{PG} J(\theta) \\ \nabla_{\theta}^{PGT} J(\theta) = G_{\theta} w \end{array}} \right\} \nabla_{\theta}^{NG} J(\theta) = F_{\theta}^{-1} G_{\theta} w$$

From PGT

$$\begin{array}{l}
 G_{\theta} = E_{p_{\theta}(\tau)} \left[ \sum_{t=0}^T \nabla_{\theta} (\log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t)) \nabla_{\theta} (\log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t))^T \right] \\
 F_{\theta} = E_{p_{\theta}(\tau)} \left[ \sum_{t=0}^T \nabla_{\theta} (\log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t)) \nabla_{\theta} (\log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t))^T \right]
 \end{array}
 \left. \vphantom{\begin{array}{l} G_{\theta} = E_{p_{\theta}(\tau)} \left[ \sum_{t=0}^T \nabla_{\theta} (\log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t)) \nabla_{\theta} (\log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t))^T \right] \\ F_{\theta} = E_{p_{\theta}(\tau)} \left[ \sum_{t=0}^T \nabla_{\theta} (\log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t)) \nabla_{\theta} (\log \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t))^T \right] \end{array}} \right\} F_{\theta} = G_{\theta}$$

$$\left. \vphantom{\begin{array}{l} \nabla_{\theta}^{NG} J(\theta) = F_{\theta}^{-1} G_{\theta} w \\ F_{\theta} = G_{\theta} \end{array}} \right\} \nabla_{\theta}^{NG} J(\theta) = w$$



# Natural Policy Gradient plus PGT with value function baseline

- Using PGT gradient with value function baseline in natural gradient format yields

$$\nabla_{\theta}^{NG} J(\theta) = w$$

- We just need to compute  $w$  through minimizing ESE

$$\min_w E_{p_{\theta}} \left[ \left( A_t^{\pi}(x_t, u_t) - w^T \nabla_{\theta} \log \pi_{\theta}(u_t | x_t) \right)^2 \right]$$

# Approximating advantage function

- To solve ESE, we should have the advantage function of the current policy.

$$\min_w E_{p_\theta} \left[ \left( A_t^\pi(x_t, u_t) - w^T \nabla_\theta \log \pi_\theta(u_t | x_t) \right)^2 \right]$$

- But we don't have the advantage function explicitly

Can we compute an estimation of advantage function?



# Advantage function

- From definition of advantage function we have

$$Q_t^\pi(x_t, u_t) = A_t^\pi(x_t, u_t) + V_t^\pi(x_t)$$

- From definition of state-action value function we have

$$Q_t^\pi(x_t, u_t) = r_t(x_t, u_t) + \int V_{t+1}^\pi(x') p(x' | x_t, u_t) dx'$$

- Combining these two formulas

$$A_t^\pi(x_t, u_t) + V_t^\pi(x_t) = r_t(x_t, u_t) + \int V_{t+1}^\pi(x') p(x' | x_t, u_t) dx'$$

# Advantage function estimation

- We got a Bellman like equation for the advantage function:

$$A_t^\pi(x_t, u_t) + V_t^\pi(x_t) = r_t(x_t, u_t) + \int V_{t+1}^\pi(x') p(x' | x_t, u_t) dx'$$

- For the samples derived from rollout, we can write

$$\tilde{A}_t^\pi(x_t, u_t) + \tilde{V}_t^\pi(x_t) = r_t(x_t, u_t) + \tilde{V}_{t+1}^\pi(x_{t+1}) + \varepsilon_t$$

For every time step  $t$

← ← ← ←

Estimated values error

Rollout: the trajectory of state and actions yields from execution of policy in the environment

$$\tau : x_0, u_0, x_1, u_1, \dots, x_t, u_t, x_{t+1}, u_{t+1}, \dots, x_{H-1}, u_{H-1}, x_H$$



# episodic Natural Actor Critic (eNAC)

- Use the estimated advantage function in each time step and sum them up

$$\begin{aligned}
 \tilde{A}_0^\pi(x_0, u_0) + \tilde{V}_0^\pi(x_0) &= r_0(x_0, u_0) + \cancel{\tilde{V}_1^\pi(x_1)} + \varepsilon_0 \\
 \tilde{A}_1^\pi(x_1, u_1) + \cancel{\tilde{V}_1^\pi(x_1)} &= r_1(x_1, u_1) + \cancel{\tilde{V}_2^\pi(x_2)} + \varepsilon_1 \\
 &\vdots \\
 \tilde{A}_t^\pi(x_t, u_t) + \cancel{\tilde{V}_t^\pi(x_t)} &= r_t(x_t, u_t) + \cancel{\tilde{V}_{t+1}^\pi(x_{t+1})} + \varepsilon_t \\
 \tilde{A}_{t+1}^\pi(x_{t+1}, u_{t+1}) + \cancel{\tilde{V}_{t+1}^\pi(x_{t+1})} &= r_{t+1}(x_{t+1}, u_{t+1}) + \cancel{\tilde{V}_{t+2}^\pi(x_{t+2})} + \varepsilon_{t+1} \\
 &\vdots \\
 \tilde{A}_{T-1}^\pi(x_T, u_T) + \cancel{\tilde{V}_T^\pi(x_T)} &= r_T(x_T) + \varepsilon_T
 \end{aligned}$$

+

---


$$\sum_{t=0}^{T-1} \tilde{A}_t^\pi(x_t, u_t) + \tilde{V}_0^\pi(x_0) = \sum_{t=0}^{T-1} r_t(x_t, u_t) + \sum_{t=0}^{T-1} \varepsilon_t$$



# Continued

- Using function approximation for advantage function

$$\sum_{t=0}^{T-1} \tilde{A}_t^\pi(x_t, u_t) + \tilde{V}_0^\pi(x_0) = \sum_{t=0}^T r_t(x_t, u_t) + \sum_{t=0}^T \varepsilon_t \quad \tilde{A}_t^\pi(x_t, u_t) \approx w^T \nabla_\theta \log \pi_\theta(u_t | x_t)$$



$$\sum_{t=0}^{T-1} w^T \nabla_\theta \log \pi_\theta(u_t | x_t) + \tilde{V}_0^\pi(x_0) = \sum_{t=0}^T r_t(x_t, u_t) + \sum_{t=0}^T \varepsilon_t$$

$$w^T \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(u_t | x_t) + \tilde{V}_0^\pi(x_0) = \sum_{t=0}^T r_t(x_t, u_t) + \sum_{t=0}^T \varepsilon_t$$



# Continued

- Now what about the value function for initial time

$$w^T \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(u_t | x_t) - \tilde{V}_0^{\pi}(x_0) = \sum_{t=0}^{T-1} r_t(x_t, u_t) + \sum_{t=0}^{T-1} \varepsilon_t$$

- We need to approximate the initial value function as well

$$\tilde{V}_0^{\pi}(x_0) \approx v^T \varphi(x_0)$$

- If the agent is always initialized in a specific state, the base function is simply one and  $v$  is accumulated reward of the trajectory
- If the agent is initialized in random states, the base function should be a function of state vector



# Continued

- Using function approximation for value function

$$w^T \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(u_t | x_t) + v^T \varphi(x_0) = \sum_{t=0}^T r_t(x_t, u_t) + \sum_{t=0}^T \varepsilon_t$$

- In the vector form we can write

$$\begin{bmatrix} w \\ v \end{bmatrix}^T \begin{bmatrix} \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(u_t | x_t) \\ \varphi(x_0) \end{bmatrix} = \sum_{t=0}^T r_t(x_t, u_t) + \sum_{t=0}^T \varepsilon_t$$

- To abbreviate the notation

$$\underbrace{\begin{bmatrix} w \\ v \end{bmatrix}}_{\text{Parameter vector}}^T \underbrace{\begin{bmatrix} \phi \\ \varphi(x_0) \end{bmatrix}}_{\text{Base function}} = R + \varepsilon$$

← Error of Estimation & Approximation  
← acc. reward of trajectory

$$\phi = \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(u_t | x_t)$$

$$R = \sum_{t=0}^T r_t(x_t, u_t)$$

$$\varepsilon = \sum_{t=0}^T \varepsilon_t$$



# eNAC Algorithm

- To reduce the error, we should use information from several rollouts (say N rollouts)

$$\begin{bmatrix} \phi^{1T} & \varphi(x_0^1)^T \end{bmatrix} \begin{bmatrix} w \\ v \end{bmatrix} = R^1 + \varepsilon^1$$

$$\vdots$$

$$\begin{bmatrix} \phi^{iT} & \varphi(x_0^i)^T \end{bmatrix} \begin{bmatrix} w \\ v \end{bmatrix} = R^i + \varepsilon^i$$

$$\vdots$$

$$\begin{bmatrix} \phi^{NT} & \varphi(x_0^N)^T \end{bmatrix} \begin{bmatrix} w \\ v \end{bmatrix} = R^N + \varepsilon^N$$

Use Least Square methods  
to estimate the  
parameter vector



# eNAC algorithm

- Using the least square method over  $N$  rollouts, we will have

$$\Psi = \begin{bmatrix} \Phi^{1^T} \\ \vdots \\ \Phi^{N^T} \end{bmatrix}, \quad \Phi^i = \begin{bmatrix} \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(u_t^i | x_t^i) \\ \varphi(x_0^i) \end{bmatrix}, \quad R = \begin{bmatrix} R^1 \\ \vdots \\ R^N \end{bmatrix}, \quad R^i = \sum_{t=0}^T r_t(x_t^i, u_t^i)$$

$$\begin{bmatrix} w \\ v \end{bmatrix} = (\Psi^T \Psi)^{-1} \Psi^T R$$

---

**Algorithm** Episodic Natural Actor Critic
 

---

Input: Policy parametrization  $\theta$ ,

$$\text{data-set } \mathcal{D} = \left\{ \mathbf{x}_{1:T}^{[i]}, \mathbf{u}_{1:T-1}^{[i]}, r_{1:T}^{[i]} \right\}_{i=1 \dots N}$$

**for** each sample  $i = 1 \dots N$  **do**

  Compute returns:  $R^{[i]} = \sum_{t=0}^T r_t^{[i]}$

$$\text{  Compute features: } \boldsymbol{\psi}^{[i]} = \begin{bmatrix} \sum_{t=0}^{T-1} \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}} \left( \mathbf{u}_t^{[i]} \mid \mathbf{x}_t^{[i]}, t \right) \\ \boldsymbol{\varphi}(\mathbf{x}_0^{[i]}) \end{bmatrix}$$

**end for**

Fit advantage function and initial value function

$$\mathbf{R} = \left[ R^{[1]}, \dots, R^{[N]} \right]^T, \quad \boldsymbol{\Psi} = \left[ \boldsymbol{\psi}^{[1]}, \dots, \boldsymbol{\psi}^{[N]} \right]^T$$

$$\begin{bmatrix} \mathbf{w} \\ \mathbf{v} \end{bmatrix} = (\boldsymbol{\Psi}^T \boldsymbol{\Psi})^{-1} \boldsymbol{\Psi}^T \mathbf{R}$$

return  $\nabla_{\boldsymbol{\theta}}^{\text{eNAC}} J_{\boldsymbol{\theta}} = \mathbf{w}$

---

# Credits

material from:

Reinforcement Learning of motor skill with policy  
gradients

A Survey on Policy Search for Robotics



# Optimal and Learning Control for Autonomous Robots Lecture 10



**A D R L**

Jonas Buchli  
Agile & Dexterous Robotics Lab



Swiss National  
Centre of Competence  
in Research



# Exercise 3

online until end of week!





# Office hour

No office hour this week!

Next week, two office hours:  
Thu 13h-15, 17:30-18:30

# Next week's lecture

## RL Recap



# Reading

There is no book!

Learning variable impedance control.

Buchli, Stulp, Theodorou, Schaal, IJRR 30(7),  
820-33



# Outline

Natural Gradient  
Natural Actor Critic

Path Integral Stochastic Optimal Control  
Policy Improvements with Path Integrals



# eNAC



# Avoid calculating explicit gradient?

Idea: Modify current best guess for optimal controls -  
Pick best seen outcome as new best guess

- This works both as global, one step algorithm
- and local, iterative algorithm
- one step algorithms run into curse of dimensionality, iterative work in practice but give local optimum
- No need for step-size! Complete update step extracted from data!

PI<sup>2</sup>

PoWER



# Stochastic optimal control



$$R(\tau_i) = \phi_{t_N} + \int_{t_i}^{t_N} r_t dt$$

$$r_t = r(\mathbf{x}_t, \mathbf{u}_t, t) = q_t + \frac{1}{2} \mathbf{u}_t^T \mathbf{R} \mathbf{u}_t$$

q arbitrary function of x,t (but not u)

state dependent input gain matrix

nonlinear system dynamics

Linear in controls and noise

$$\dot{\mathbf{x}}_t = \mathbf{f}(\mathbf{x}_t, t) + \mathbf{G}(\mathbf{x}_t) (\mathbf{u}_t + \boldsymbol{\varepsilon}_t)$$

Noise: meanfree, gaussian

$$\mathbf{f}_t + \mathbf{G}_t (\mathbf{u}_t + \boldsymbol{\varepsilon}_t)$$

$$V(x_t) = \min_{u_t} E_{\tau} [R(\tau)]$$

$$u^*(x, t) = \operatorname{argmin}_{u_t} E_{\tau} [R(\tau)]$$



## Lecture 3: LQR

quadratic cost

$$V = \frac{1}{2} \Delta \mathbf{x}^T(t_f) \phi_{\mathbf{xx}}(t_f) \Delta \mathbf{x}(t_f)$$

$$+ \frac{1}{2} \int_{t_0}^{t_f} \left\{ [\Delta \mathbf{x}^T(t) \ \Delta \mathbf{u}^T(t)] \begin{bmatrix} \mathbf{Q}(t) & \mathbf{M}(t) \\ \mathbf{M}^T(t) & \mathbf{R}(t) \end{bmatrix} \begin{bmatrix} \Delta \mathbf{x}(t) \\ \Delta \mathbf{u}(t) \end{bmatrix} \right\}$$

linear dynamics

$$\Delta \dot{\mathbf{x}}(t) = \mathbf{F}(t) \Delta \mathbf{x}(t) + \mathbf{G}(t) \Delta \mathbf{u}(t),$$



# HJB equation



Buchli - OLCAR - 2013



# Solving Quadratic HJB Equation

LHS

$$\frac{\partial V^*}{\partial t} [\Delta \mathbf{x}^*(t), t]$$

①

~~$$\frac{1}{2} \Delta \mathbf{x}^{*T}(t) \dot{\mathbf{P}}(t) \Delta \mathbf{x}^*(t)$$~~

RHS

$$-\min_{\mathbf{u}} \mathcal{H} [\Delta \mathbf{x}^*(t), \Delta \mathbf{u}(t), t]$$

②

$$\frac{\partial}{\partial \mathbf{u}} (\cdot) = 0$$

③

④

 $\mathbf{u}^*$ 

⑤

~~$$\frac{\partial V^*}{\partial \Delta \mathbf{x}} [\Delta \mathbf{x}^*(t), t] = \Delta \mathbf{x}^{*T}(t) \mathbf{P}(t)$$~~

# Find Value Function and Optimal Controls

Find right side of HJB:

- ① by partial derivative with respect to controls and setting it to 0
- ~~② Use partial derivative of Quadratic Ansatz to substitute partial derivative of  $V$  in respect to  $x$~~
- ③ solve in this expression for  $u$ : yields  $u^*$  (optimal control!)
- ~~④ substitute  $u^*$  back into HJB, solve for unknown matrix  $P$~~

$$V(x_t) = \min_{u_t} E_\tau [R(\tau)]$$

$$\int p(\tau) R(\tau) d\tau \quad \int p(\tau) \left( \frac{1}{\lambda} \phi + \frac{1}{\lambda} \int r dt \right) d\tau$$

Discretize and use EM-like idea: PoWER -  
 Problem: pseudo-probability - restriction on  
 cost function

Other idea: Treat probability as a diffusion  
 process - Connection with statistical physics  
 Forward dynamics! Sampling (Monte Carlo)!



# Derivation of stochastic HJB

# Stochastic principle of optimality

Principle of optimality

$$V^*(t_1) = E \left\{ \phi[\mathbf{x}^*(t_f), t_f] - \int_{t_f}^{t_1} \mathcal{L}[\mathbf{x}^*(t), \mathbf{u}^*(t), t] dt \right\}$$

Total time derivative:

$$\frac{dV^*(t_1)}{dt} = - E \{ \mathcal{L}[\mathbf{x}^*(t_1), \mathbf{u}^*(t_1), t_1] \}$$

Measurements are deterministic

$$\frac{dV^*(t_1)}{dt} = - \mathcal{L}[\mathbf{x}^*(t_1), \mathbf{u}^*(t_1), t_1]$$

[St] p 422/23

## Can also write total time derivative as Taylor Series

$$\begin{aligned}\frac{dV^*}{dt} \Delta t &= E \left\{ \frac{\partial V^*}{\partial t} \Delta t + \frac{\partial V^*}{\partial \mathbf{x}} \dot{\mathbf{x}} \Delta t + \frac{1}{2} \left[ \dot{\mathbf{x}}^T \frac{\partial^2 V^*}{\partial \mathbf{x}^2} \dot{\mathbf{x}} \right] \Delta t^2 + \dots \right\} \\ &= E \left[ V_t^* \Delta t + V_{\mathbf{x}}^* (\mathbf{f} + \mathbf{L}\mathbf{w}) \Delta t + \frac{1}{2} (\mathbf{f} + \mathbf{L}\mathbf{w})^T V_{\mathbf{xx}}^* (\mathbf{f} + \mathbf{L}\mathbf{w}) \Delta t^2 \right]\end{aligned}$$

Functions of  $\mathbf{x}(t)$  equal their own expectations, and  $E[\mathbf{w}(t)] = \mathbf{0}$ . Dividing by  $\Delta t$ , and replacing the third term by its trace, the time derivative is

$$\begin{aligned}\frac{dV^*}{dt} &= V_t^* + V_{\mathbf{x}}^* \mathbf{f} + \frac{1}{2} \text{Tr} \{ E [ (\mathbf{f} + \mathbf{L}\mathbf{w})^T V_{\mathbf{xx}}^* (\mathbf{f} + \mathbf{L}\mathbf{w}) ] \Delta t \} \\ &= V_t^* + V_{\mathbf{x}}^* \mathbf{f} + \frac{1}{2} \text{Tr} \{ E [ V_{\mathbf{xx}}^* (\mathbf{f} + \mathbf{L}\mathbf{w}) (\mathbf{f} + \mathbf{L}\mathbf{w})^T ] \Delta t \}\end{aligned}$$

# Stochastic HJB

$\mathbf{f}$ ,  $\mathbf{w}$  uncorrelated

$$\begin{aligned}\frac{dV^*}{dt} &= V_t^* + V_x^* \mathbf{f} + \frac{1}{2} \lim_{\Delta t \rightarrow 0} \text{Tr}\{V_{xx}^* [E(\mathbf{f}\mathbf{f}^T) \Delta t + \mathbf{L}E(\mathbf{w}\mathbf{w}^T)\mathbf{L}^T \Delta t]\} \\ &= V_t^* + V_x^* \mathbf{f} + \frac{1}{2} \text{Tr}(V_{xx}^* \mathbf{L}\mathbf{W}\mathbf{L}^T)\end{aligned}$$

plug in and rearrange

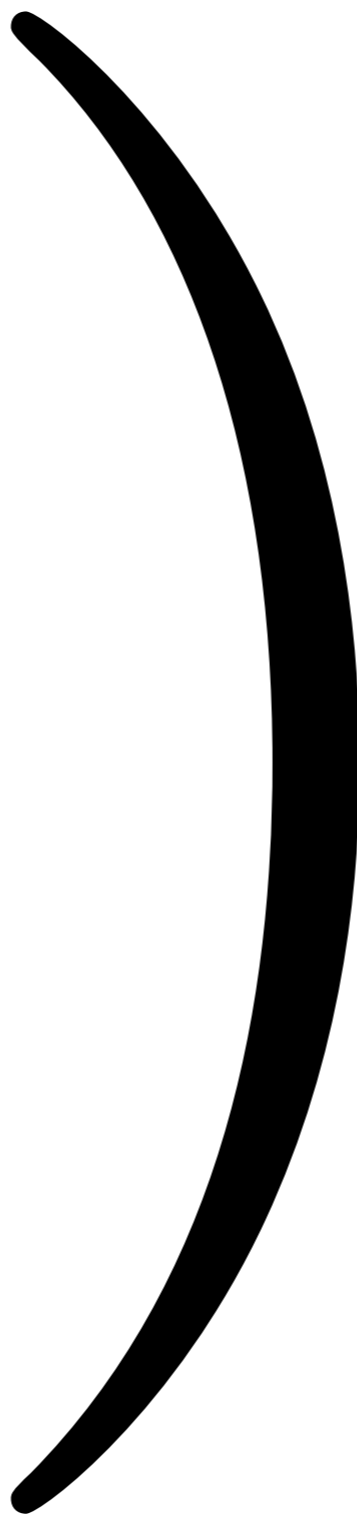
$$\begin{aligned}V_t^*(t) &= -\min_{\mathbf{u}} \{ \mathcal{L}[(\mathbf{x}^*(t), \mathbf{u}(t), t)] + V_x^* \mathbf{f}[\mathbf{x}^*(t), \mathbf{u}(t), t] \\ &\quad + \frac{1}{2} \text{Tr}[V_{xx}^* \mathbf{L}(t)\mathbf{W}(t)\mathbf{L}^T(t)] \}\end{aligned}$$

Terminal condition:

---

starting value for evaluation of  $V^*(t)$  is  $E\{\phi[\mathbf{x}(t_f), t_f]\}$ , which is  $\phi[\mathbf{x}(t_f), t_f]$  because  $\mathbf{x}(t_f)$  can be measured without error.





Buchli - OLCAR - 2013



# nonlinear HJB

$$-\partial_t V_t = \min_{\mathbf{u}} \left( r_t + (\nabla_{\mathbf{x}} V_t)^T \mathbf{F}_t + \frac{1}{2} \text{trace} \left( (\nabla_{\mathbf{xx}} V_t) \mathbf{G}_t \Sigma_{\varepsilon} \mathbf{G}_t^T \right) \right)$$

$$\mathbf{F}_t = \mathbf{f}(\mathbf{x}_t, t) + \mathbf{G}(\mathbf{x}_t) \mathbf{u}_t$$

① gradient of RHS = 0 yields

③

$$\mathbf{u}(\mathbf{x}_t) = \mathbf{u}_t = -\mathbf{R}^{-1} \mathbf{G}_t^T (\nabla_{x_t} V_t)$$

④ substitute opt. control back into HJB  $\Rightarrow$

$$-\partial_t V_t = q_t + (\nabla_{\mathbf{x}} V_t)^T \mathbf{f}_t - \frac{1}{2} (\nabla_{\mathbf{x}} V_t)^T \mathbf{G}_t \mathbf{R}^{-1} \mathbf{G}_t^T (\nabla_{\mathbf{x}} V_t) + \frac{1}{2} \text{trace} \left( (\nabla_{\mathbf{xx}} V_t) \mathbf{G}_t \Sigma_{\varepsilon} \mathbf{G}_t^T \right)$$

**Nonlinear PDE!**



# log transform

$$-\partial_t V_t = q_t + (\nabla_{\mathbf{x}} V_t)^T \mathbf{f}_t - \frac{1}{2} (\nabla_{\mathbf{x}} V_t)^T \mathbf{G}_t \mathbf{R}^{-1} \mathbf{G}_t^T (\nabla_{\mathbf{x}} V_t) + \frac{1}{2} \text{trace} \left( (\nabla_{\mathbf{xx}} V_t) \mathbf{G}_t \Sigma_{\varepsilon} \mathbf{G}_t^T \right)$$

Nonlinear PDE!

$$V_t = -\lambda \log \Psi_t$$

$\Rightarrow$

$$\partial_t V_t = -\lambda \frac{1}{\Psi_t} \partial_t \Psi_t,$$

$$\nabla_{\mathbf{x}} V_t = -\lambda \frac{1}{\Psi_t} \nabla_{\mathbf{x}} \Psi_t,$$

$$\nabla_{\mathbf{xx}} V_t = \lambda \frac{1}{\Psi_t^2} \nabla_{\mathbf{x}} \Psi_t \nabla_{\mathbf{x}} \Psi_t^T - \lambda \frac{1}{\Psi_t} \nabla_{\mathbf{xx}} \Psi_t$$

$$\frac{\lambda}{\Psi_t} \partial_t \Psi_t = q_t - \frac{\lambda}{\Psi_t} (\nabla_{\mathbf{x}} \Psi_t)^T \mathbf{f}_t - \frac{\lambda^2}{2\Psi_t^2} (\nabla_{\mathbf{x}} \Psi_t)^T \mathbf{G}_t \mathbf{R}^{-1} \mathbf{G}_t^T (\nabla_{\mathbf{x}} \Psi_t) + \frac{1}{2} \text{trace} (\Gamma)$$

$$\Gamma = \left( \lambda \frac{1}{\Psi_t^2} \nabla_{\mathbf{x}} \Psi_t \nabla_{\mathbf{x}} \Psi_t^T - \lambda \frac{1}{\Psi_t} \nabla_{\mathbf{xx}} \Psi_t \right) \mathbf{G}_t \Sigma_{\varepsilon} \mathbf{G}_t^T$$

# structure of control cost linked to noise

$$\frac{\lambda}{\Psi_t} \partial_t \Psi_t = q_t - \frac{\lambda}{\Psi_t} (\nabla_{\mathbf{x}} \Psi_t)^T \mathbf{f}_t - \frac{\lambda^2}{2\Psi_t^2} (\nabla_{\mathbf{x}} \Psi_t)^T \mathbf{G}_t \mathbf{R}^{-1} \mathbf{G}_t^T (\nabla_{\mathbf{x}} \Psi_t) + \frac{1}{2} \text{trace}(\Gamma)$$

$$\Gamma = \left( \lambda \frac{1}{\Psi_t^2} \nabla_{\mathbf{x}} \Psi_t \nabla_{\mathbf{x}} \Psi_t^T - \lambda \frac{1}{\Psi_t} \nabla_{\mathbf{xx}} \Psi_t \right) \mathbf{G}_t \Sigma_{\varepsilon} \mathbf{G}_t^T$$

$$\text{trace}(\Gamma) = \lambda \frac{1}{\Psi_t^2} \text{trace}(\nabla_{\mathbf{x}} \Psi_t^T \mathbf{G}_t \Sigma_{\varepsilon} \mathbf{G}_t \nabla_{\mathbf{x}} \Psi_t) - \lambda \frac{1}{\Psi_t} \text{trace}(\nabla_{\mathbf{xx}} \Psi_t \mathbf{G}_t \Sigma_{\varepsilon} \mathbf{G}_t^T)$$

$$\lambda \mathbf{R}^{-1} = \Sigma_{\varepsilon}$$

$$\lambda \mathbf{G}_t \mathbf{R}^{-1} \mathbf{G}_t^T = \mathbf{G}_t \Sigma_{\varepsilon} \mathbf{G}_t^T = \Sigma(\mathbf{x}_t) = \Sigma_t$$

$$-\partial_t \Psi_t = -\frac{1}{\lambda} q_t \Psi_t + \mathbf{f}_t^T (\nabla_{\mathbf{x}} \Psi_t) + \frac{1}{2} \text{trace}((\nabla_{\mathbf{xx}} \Psi_t) \mathbf{G}_t \Sigma_{\varepsilon} \mathbf{G}_t^T)$$



linear!



# linear HJB

$$-\partial_t \Psi_t = -\frac{1}{\lambda} q_t \Psi_t + \mathbf{f}_t^T (\nabla_{\mathbf{x}} \Psi_t) + \frac{1}{2} \text{trace} \left( (\nabla_{\mathbf{xx}} \Psi_t) \mathbf{G}_t \Sigma_{\epsilon} \mathbf{G}_t^T \right)$$

linear, but still no analytic solution for arbitrary  $q(\mathbf{x}, t)$

solve **backward**

terminal condition :  $\Psi_{t_N} = \exp \left( -\frac{1}{\lambda} \phi_{t_N} \right)$

**Feynman-Kac Theorem:** Can write solution of PDE as  
Expectation over stochastic forward dynamics

$$\Psi_{t_i} = E_{\tau_i} \left( \Psi_{t_N} e^{-\int_{t_i}^{t_N} \frac{1}{\lambda} q_t dt} \right) = E_{\tau_i} \left[ \exp \left( -\frac{1}{\lambda} \phi_{t_N} - \frac{1}{\lambda} \int_{t_i}^{t_N} q_t dt \right) \right]$$

forward! ... but stochastic

Remember the forward search in the discrete  
state, discrete time problem (Lect. 2)



# Expectations over paths

$$\Psi_{t_i} = E_{\tau_i} \left( \Psi_{t_N} e^{-\int_{t_i}^{t_N} \frac{1}{\lambda} q_t dt} \right) = E_{\tau_i} \left[ \exp \left( -\frac{1}{\lambda} \phi_{t_N} - \frac{1}{\lambda} \int_{t_i}^{t_N} q_t dt \right) \right]$$

forward!

... but stochastic

$$\int p(\tau) \exp \left( -\frac{1}{\lambda} \phi - \frac{1}{\lambda} \int q dt \right) d\tau$$

$$\tau = x(t \dots t_N) \sim p(x, u)$$

an instance of a random path segment (a random 'number', but in spaces of functions)

$$E[X] = \int x p(x) dx$$

$$x = f(t)$$

Continuous time,  $x$  is function of time



# Major difficulty: Definition of stochastic processes in continuous time!

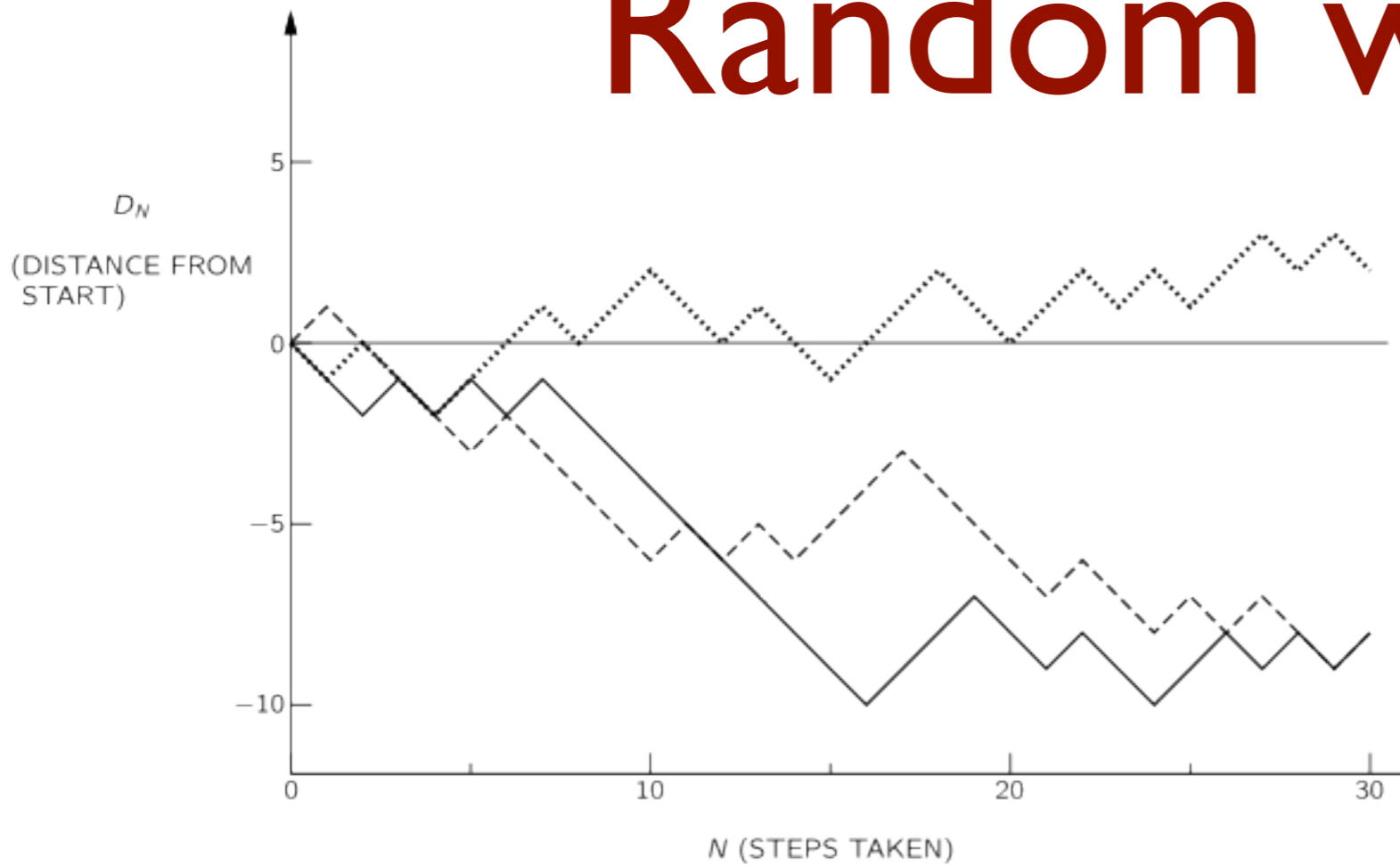
		states	
		discrete	continuous
time	discrete		$\begin{array}{c} dx \rightarrow 0 \\ \leftarrow \\ \int \rightarrow \Sigma \end{array}$
	continuous		$\int \rightarrow \Sigma \begin{array}{c} \uparrow \\ \downarrow \end{array} dt \rightarrow 0$ <b>Lots of gnarly math!</b>

# Continuous Random Processes

$$dx = f(x, u) dt + F(x, u) d\omega$$

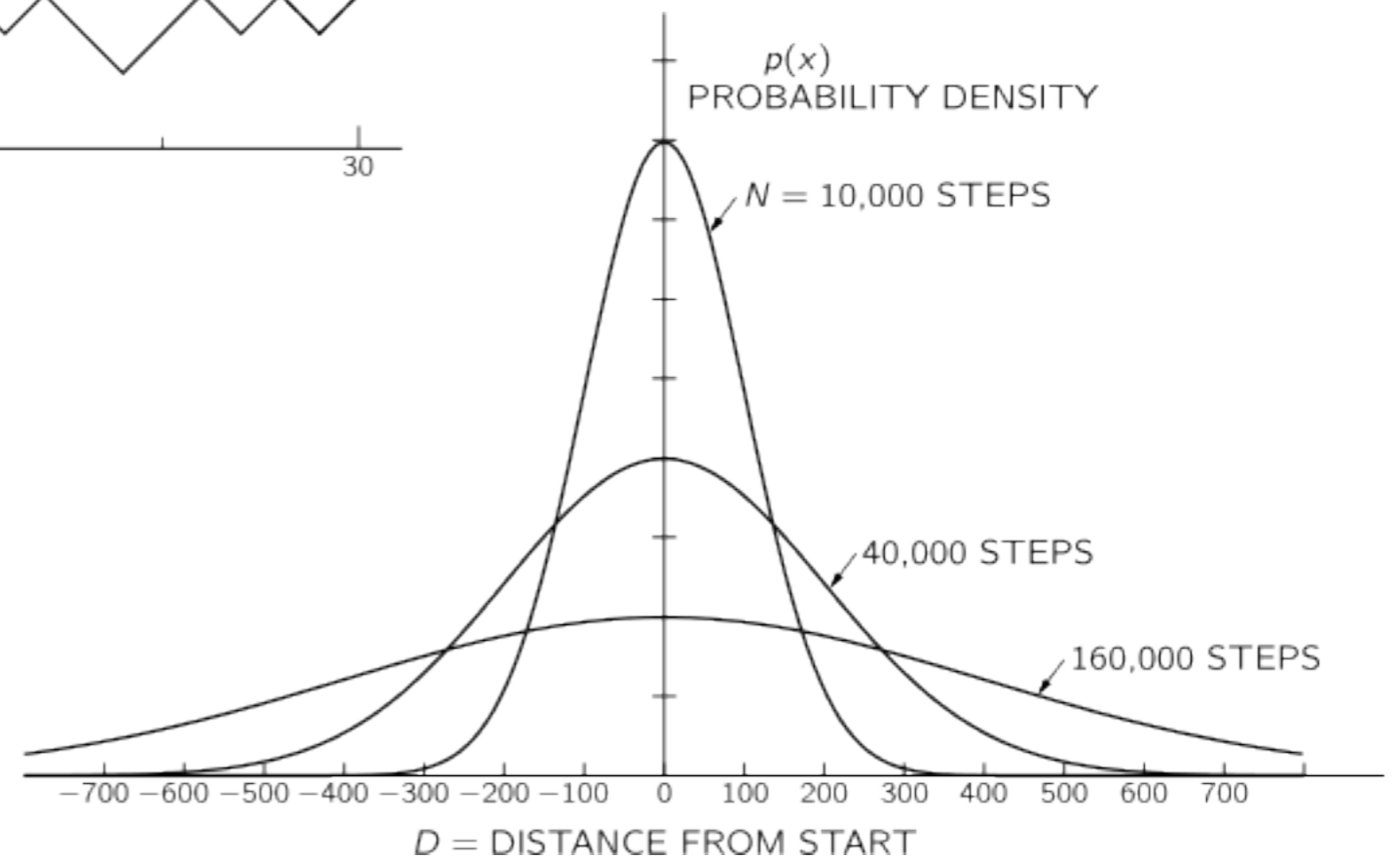


# Random walks



‘paths diffuse over time’

‘density’



# Continuous decision processes

Take random walk and take limits

$dx \rightarrow 0$       probability densities

$dt \rightarrow 0$       probability flow

for all times  $\int p(x) dx = 1$       conservation law!

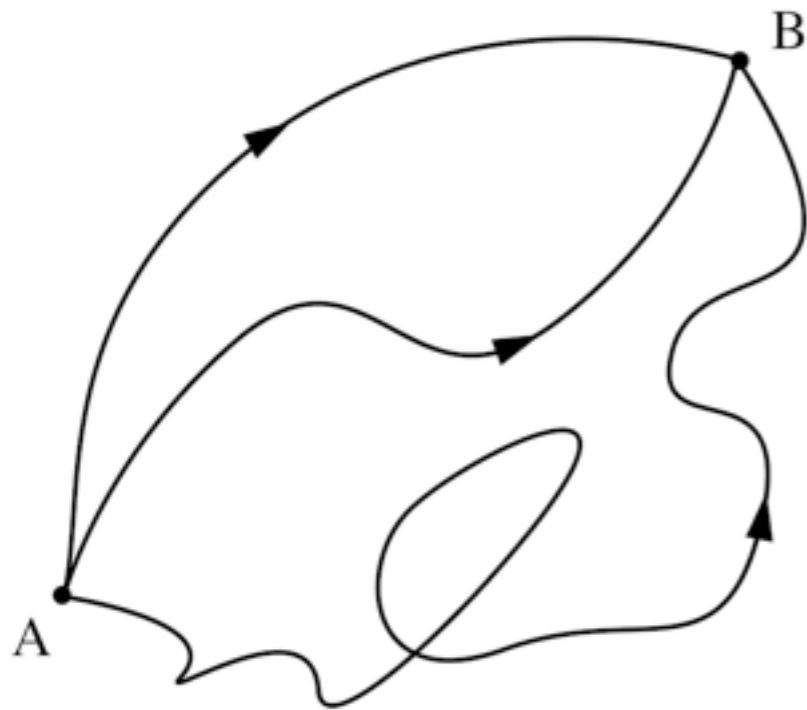
Conserved flow?

You know how to do that!

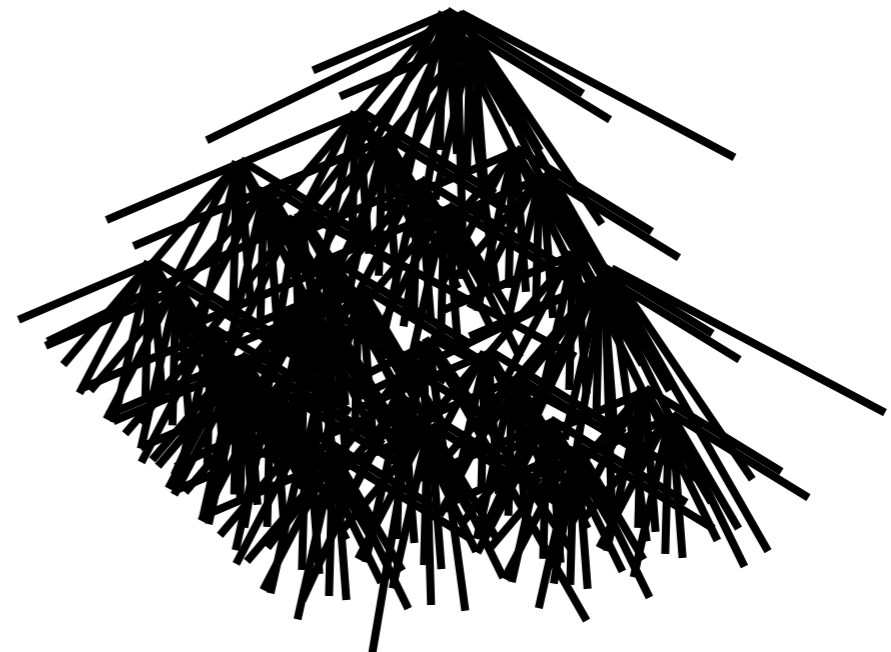


# Comparison to graphs

can think of all possibilities of a random walk as graph



There are several ways to end up in a certain state, each path has an associated probability



When does 'branching' occur?

Idea: do discrete time and take limit



$$dt \rightarrow 0$$

# Probabilistic Dynamics

Discrete time:

Markov chains

Master Equation

Continuous time:

Jumps: Continuous-time Markov chain

Smooth: Markov Process

Fokker-Planck

cf. (Heat) Diffusion



# Fokker-Planck Equation

(the most interesting equation in the world?)

PDE for time evolution of probability distribution

$$\frac{\partial}{\partial t} p(x, t) = - \frac{\partial}{\partial x} [\mu(x, t) p(x, t)] + \frac{\partial^2}{\partial x^2} [D(x, t) p(x, t)]$$

**Drift**

**Diffusion**

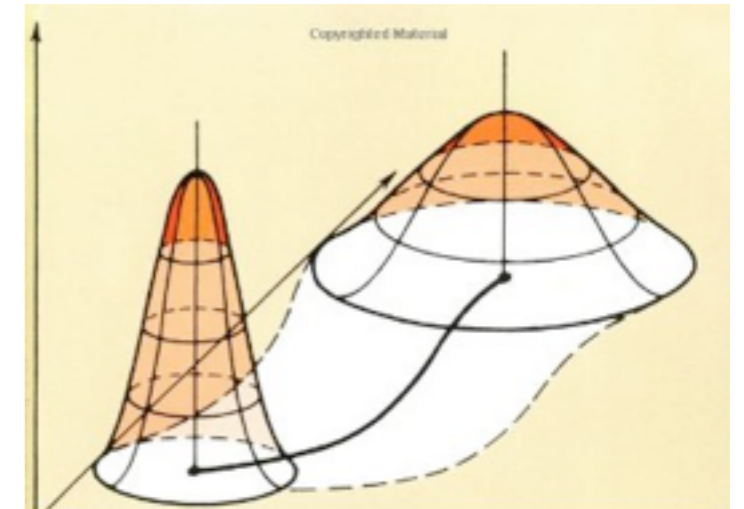
$$dX_t = dW_t. \text{ brownian motion, no drift}$$

$$\frac{\partial p(x, t)}{\partial x} = \frac{1}{2} \frac{\partial^2 p(x, t)}{\partial x^2}$$

$$\Rightarrow p(x, t) = \frac{1}{\sqrt{2\pi t}} e^{-\frac{x^2}{2t}}$$

conservation law!

$$\int p(x) dx = 1$$



cf. Fluid Dynamics  
Heat and Charge diffusion  
cf. Particle filters

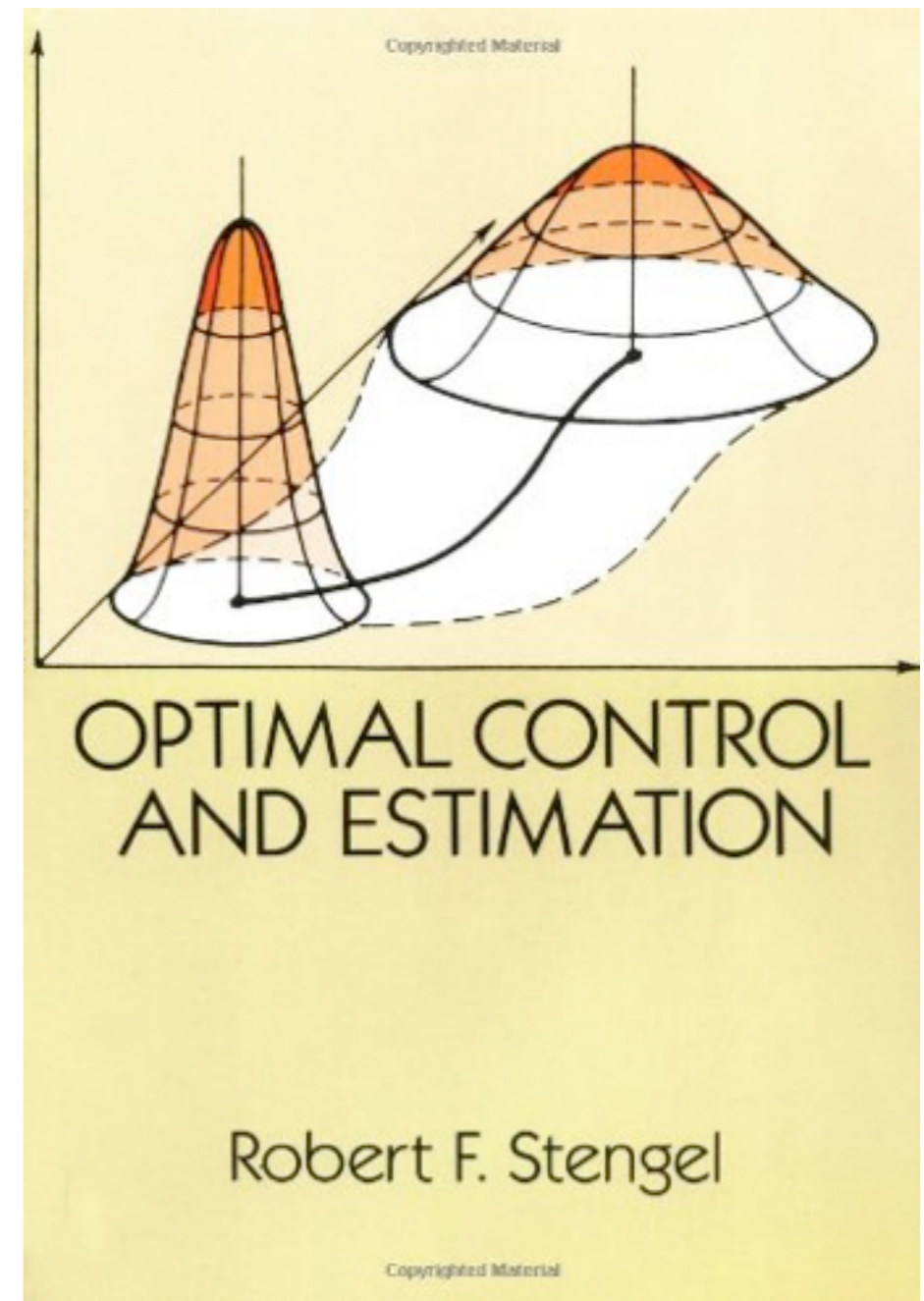
Finance, Biology, Chemistry, Physics, Sociology,  
Anthropology, Control & Machine Learning



Zürich

# Stochastic Control

## 'Controlled Diffusion' Controlled Brownian Motion



# Example

High dimensional continuous state actions spaces with stochastic dynamics

Optimal(?) control in fluids

Approach: Computational Fluid dynamics & Evolutionary Algorithm









$C = \text{'don't eat me!'}$

$$\frac{\partial C}{\partial \theta} = \frac{\partial \text{'don't eat me!'}}{\partial \text{'how to flap???'}}$$



$$\rho \left( \frac{\partial \mathbf{v}}{\partial t} + \mathbf{v} \cdot \nabla \mathbf{v} \right) = -\nabla p + \nabla \cdot \mathbf{T} + \mathbf{f}$$



$$V_t = -\lambda \log \Psi_t$$

$$\Psi_{t_i} = E_{\tau_i} \left( \Psi_{t_N} e^{-\int_{t_i}^{t_N} \frac{1}{\lambda} q_t dt} \right) = E_{\tau_i} \left[ \exp \left( -\frac{1}{\lambda} \phi_{t_N} - \frac{1}{\lambda} \int_{t_i}^{t_N} q_t dt \right) \right]$$

forward! ... but stochastic

$$\int p(\tau) \exp \left( -\frac{1}{\lambda} \phi - \frac{1}{\lambda} \int q dt \right) d\tau$$

discretize

$$\tau_i = (\mathbf{x}_{t_i}, \mathbf{x}_{t_{i+1}}, \dots, \mathbf{x}_{t_N}) \quad d\tau_i = (d\mathbf{x}_{t_i}, \dots, d\mathbf{x}_{t_N})$$

path starting at  $t_i$  to end of episode

‘how much does it cost?’

$$\Psi_{t_i} = \lim_{dt \rightarrow 0} \int p(\tau_i | \mathbf{x}_i) \exp \left[ -\frac{1}{\lambda} \left( \phi_{t_N} + \sum_{j=i}^{N-1} q_{t_j} dt \right) \right] d\tau_i$$

‘where do i end up next?’

integrate (‘sum’) over all possible paths: ‘path integral’

$$p(\tau_i | \mathbf{x}_i) ???$$



$$\begin{aligned}
p(\boldsymbol{\tau}_i | \mathbf{x}_{t_i}) &= p(\boldsymbol{\tau}_{i+1} | \mathbf{x}_{t_i}) \\
&= p(\mathbf{x}_{t_N}, \dots, \mathbf{x}_{t_{i+1}} | \mathbf{x}_{t_i}) \\
&= \prod_{j=i}^{N-1} p(\mathbf{x}_{t_{j+1}} | \mathbf{x}_{t_j}),
\end{aligned}$$

Gaussian noise leads to

$$p(\mathbf{x}_{t_{j+1}}^{(c)} | \mathbf{x}_{t_j}) = \frac{1}{((2\pi)^l \cdot |\Sigma_{t_j}|)^{1/2}} \exp\left(-\frac{1}{2} \left\| \mathbf{x}_{t_{j+1}}^{(c)} - \mathbf{x}_{t_j}^{(c)} - \mathbf{f}_{t_j}^{(c)} dt \right\|_{\Sigma_{t_j}^{-1}}^2\right)$$

‘deviation from deterministic dynamics’

$$\begin{aligned}
\Psi_{t_i} &= \lim_{dt \rightarrow 0} \int \exp\left(-\frac{1}{\lambda} S(\boldsymbol{\tau}_i) - \log D(\boldsymbol{\tau}_i)\right) d\boldsymbol{\tau}_i^{(c)} \\
&= \lim_{dt \rightarrow 0} \int \exp\left(-\frac{1}{\lambda} Z(\boldsymbol{\tau}_i)\right) d\boldsymbol{\tau}_i^{(c)},
\end{aligned}$$

‘effect of noise variance’

$$S(\boldsymbol{\tau}_i) = \phi_{t_N} + \sum_{j=i}^{N-1} q_{t_j} dt + \frac{1}{2} \sum_{j=i}^{N-1} \left\| \frac{\mathbf{x}_{t_{j+1}}^{(c)} - \mathbf{x}_{t_j}^{(c)}}{dt} - \mathbf{f}_{t_j}^{(c)} \right\|_{\mathbf{H}_{t_j}^{-1}}^2 dt$$

$$D(\boldsymbol{\tau}_i) = \prod_{j=i}^{N-1} \left( (2\pi)^{l/2} |\Sigma_{t_j}|^{1/2} \right) \text{ ‘normalization with noise variance’}$$

# Illustration

[Psi at t\_i]

$$V_t = -\lambda \log \Psi_t$$

$$\partial_t V_t = -\lambda \frac{1}{\Psi_t} \partial_t \Psi_t,$$

$$\nabla_{\mathbf{x}} V_t = -\lambda \frac{1}{\Psi_t} \nabla_{\mathbf{x}} \Psi_t,$$

$$\nabla_{\mathbf{xx}} V_t = \lambda \frac{1}{\Psi_t^2} \nabla_{\mathbf{x}} \Psi_t \nabla_{\mathbf{x}} \Psi_t^T - \lambda \frac{1}{\Psi_t} \nabla_{\mathbf{xx}} \Psi_t$$

$$\mathbf{u}_{t_i} = -\mathbf{R}^{-1} \mathbf{G}_{t_i}^T (\nabla_{\mathbf{x}_{t_i}} V_{t_i})$$

$$\Rightarrow \mathbf{u}_{t_i} = \lambda \mathbf{R}^{-1} \mathbf{G}_{t_i} \frac{\nabla_{\mathbf{x}_{t_i}} \Psi_{t_i}}{\Psi_{t_i}}$$

Plug in  $\Psi$

$$\mathbf{u}_{t_i} = \lim_{dt \rightarrow 0} \left( \lambda \mathbf{R}^{-1} \mathbf{G}_{t_i}^T \frac{\nabla_{\mathbf{x}_{t_i}^{(c)}} \left( \int e^{-\frac{1}{\lambda} \tilde{S}(\boldsymbol{\tau}_i)} d\boldsymbol{\tau}_i^{(c)} \right)}{\int e^{-\frac{1}{\lambda} \tilde{S}(\boldsymbol{\tau}_i)} d\boldsymbol{\tau}_i^{(c)}} \right)$$

$$\mathbf{u}_{t_i} = \int P(\boldsymbol{\tau}_i) \mathbf{u}_L(\boldsymbol{\tau}_i) d\boldsymbol{\tau}_i^{(c)}$$

$$\mathbf{u}_L(\boldsymbol{\tau}_i) = -\mathbf{R}^{-1} \mathbf{G}_{t_i}^{(c)T} \lim_{dt \rightarrow 0} \left( \nabla_{\mathbf{x}_{t_i}^{(c)}} \tilde{S}(\boldsymbol{\tau}_i) \right)$$

$$P(\boldsymbol{\tau}_i) = \frac{e^{-\frac{1}{\lambda} \tilde{S}(\boldsymbol{\tau}_i)}}{\int e^{-\frac{1}{\lambda} \tilde{S}(\boldsymbol{\tau}_i)} d\boldsymbol{\tau}_i}$$

$$\mathbf{u}_L(\boldsymbol{\tau}_i) = \mathbf{R}^{-1} \mathbf{G}_{t_i}^{(c)T} \left( \mathbf{G}_{t_i}^{(c)} \mathbf{R}^{-1} \mathbf{G}_{t_i}^{(c)T} \right)^{-1} \mathbf{G}_{t_i}^{(c)} \boldsymbol{\varepsilon}_{t_i}$$

'project noise in range space of control gain'

- weighted with control cost R

Buchli - OLCAR - 2013



# Example: Naive sampling

$$M(\theta) \cdot \ddot{\theta} + C(\theta, \dot{\theta}) = \tau$$

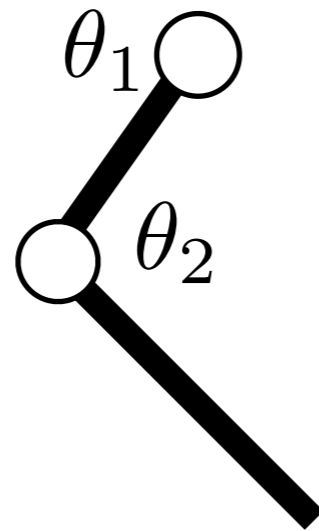
$$\ddot{\theta} = M(\theta)^{-1} \cdot (-C(\theta, \dot{\theta}) + \tau)$$

$$M(\theta) = \begin{pmatrix} d_1 + 2d_2 \cos(\theta_2) & d_3 + d_2 \cos(\theta_2) \\ d_3 + d_2 \cos(\theta_2) & d_3 \end{pmatrix} \quad (37)$$

$$C(\dot{\theta}, \theta) = \begin{pmatrix} -\dot{\theta}_2 (2\dot{\theta}_1 + \dot{\theta}_2) \\ \dot{\theta}_1^2 \end{pmatrix} \begin{matrix} d_2 \sin(\theta_2) \\ 0 \end{matrix} \quad (38)$$

$$d_1 = I_1 + I_2 + m_2 l_1^2, \quad d_2 = m_2 l_1 s_2, \quad d_3 = I_2 \quad (39)$$

Symbol	Value	Unit
$m_1$	1.4	Kg
$m_2$	1	Kg
$s_1$	0.11	m
$s_2$	0.16	m
$I_1$	0.3	Kg m <sup>2</sup>
$I_2$	0.33	Kg m <sup>2</sup>
$l_1$	0.025	m
$l_2$	0.045	m

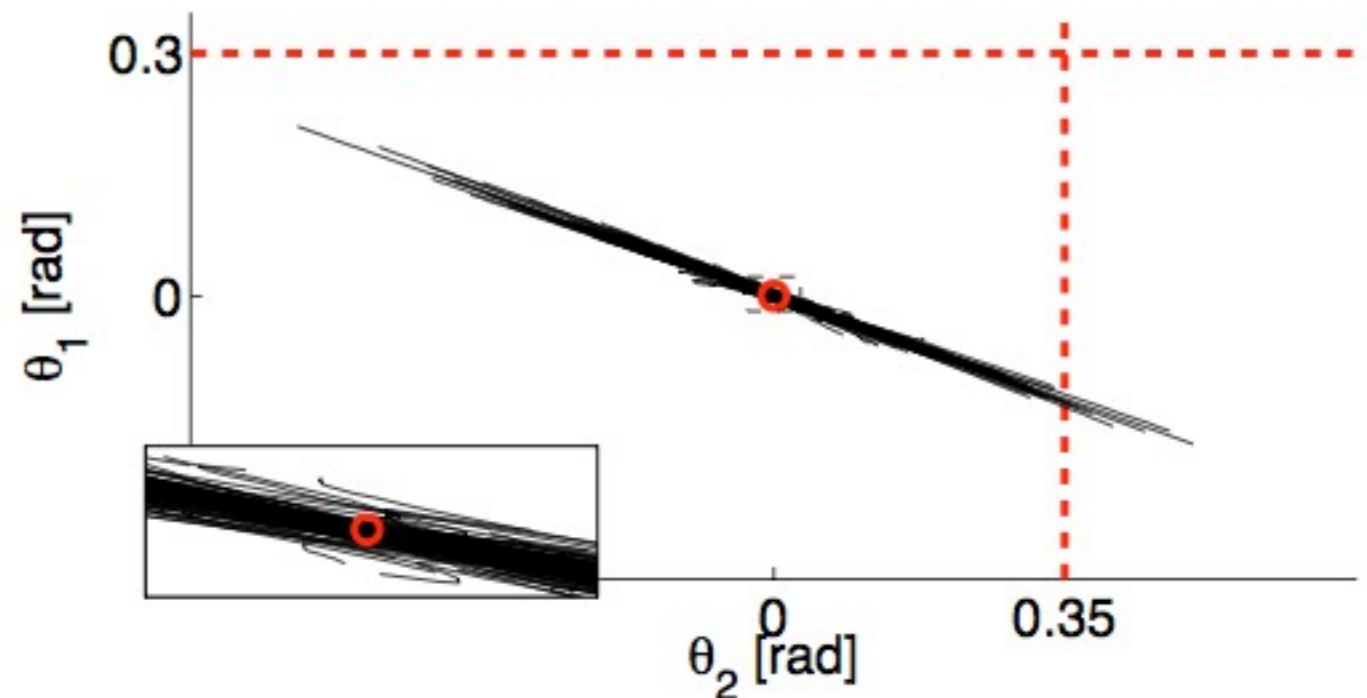


$$\dot{x} = \Phi(x) + G(x) \cdot \tau$$

$$x = \begin{pmatrix} \theta_1 & \theta_2 & \dot{\theta}_1 & \dot{\theta}_2 \end{pmatrix}^T \quad \tau = \begin{pmatrix} \tau_1 & \tau_2 \end{pmatrix}$$

$$\Phi(x) = \begin{pmatrix} \dot{\theta}_1 \\ \dot{\theta}_2 \\ -M(\theta)^{-1} \cdot C(\theta, \dot{\theta}) \end{pmatrix}, \quad G(x) = \begin{pmatrix} O_{2 \times 2} \\ M(\theta)^{-1} \end{pmatrix}$$

Path integral SOC  
requires sampling of  
passive dynamics with  
gaussian mean-free noise





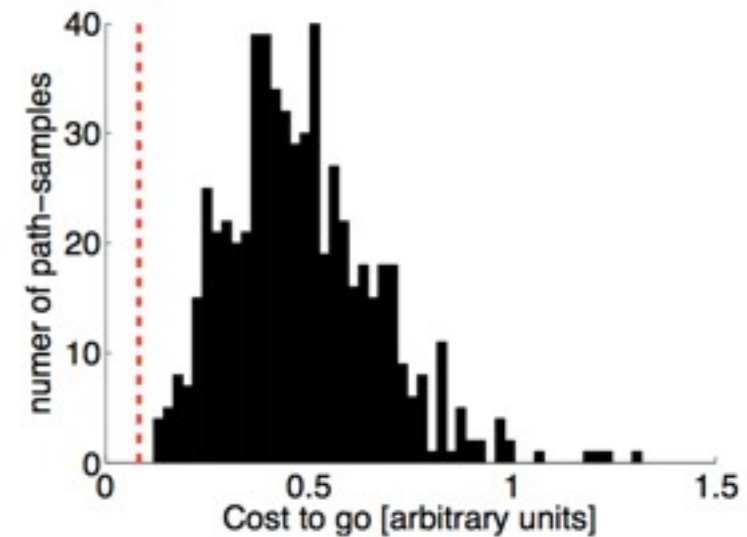
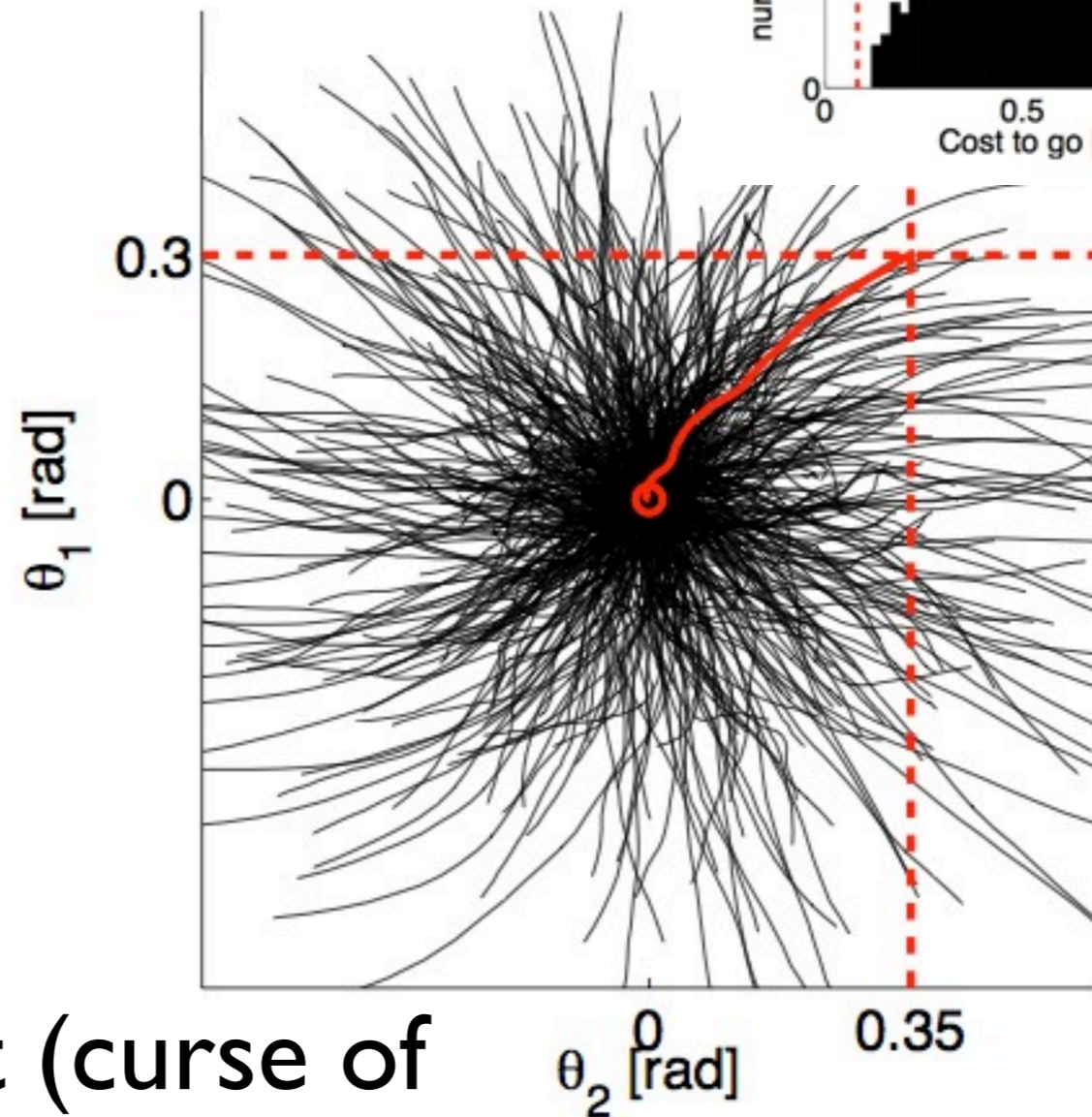
# Improved sampling...?

$$\tau_u^i = M(\theta) \cdot (\alpha_i + \epsilon_i) + C(\theta, \dot{\theta})$$

$$M(\theta) \cdot \ddot{\theta} + C(\theta, \dot{\theta}) = \tau_u$$

Sample in  
acceleration space,  
use inverse dynamics  
controllers to find  
torques:

... still not very efficient (curse of  
dimensionality still strikes, needle in a  
haystack!)



# From Model based to Model-free PISOC

$$\mathbf{u}_{t_i} = \int P(\tau_i) \mathbf{u}_L(\tau_i) d\tau_i^{(c)}$$

$$P(\tau_i) = \frac{e^{-\frac{1}{\lambda} \tilde{s}(\tau_i)}}{\int e^{-\frac{1}{\lambda} \tilde{s}(\tau_i)} d\tau_i}$$

$$\mathbf{u}_L(\tau_i) = \mathbf{R}^{-1} \mathbf{G}_{t_i}^{(c)T} \left( \mathbf{G}_{t_i}^{(c)} \mathbf{R}^{-1} \mathbf{G}_{t_i}^{(c)T} \right)^{-1} \mathbf{G}_{t_i}^{(c)} \boldsymbol{\varepsilon}_{t_i}$$

Optimal control need model: Input Gain matrix

Good sampling: need model, input gain matrix

# Policy improvements with Path Integrals - PI2

- 1) Local sampling, iterative method
- 2) Use an 'intermediate system' with known input gain matrix: parametrized policies!

$$\mathbf{a}_{t_i} = \mathbf{g}_{t_i}^T (\boldsymbol{\theta} + \boldsymbol{\varepsilon}_{t_i})$$

parameters

basis functions

noise

# Path integral SOC with parameterized policy

$$\mathbf{a}_{t_i} = \mathbf{g}_{t_i}^T (\boldsymbol{\theta} + \boldsymbol{\varepsilon}_{t_i})$$

$$\mathbf{u}_{t_i} = \int P(\boldsymbol{\tau}_i) \mathbf{u}_L(\boldsymbol{\tau}_i) d\boldsymbol{\tau}_i^{(c)}$$

$$P(\boldsymbol{\tau}_i) = \frac{e^{-\frac{1}{\lambda} \tilde{S}(\boldsymbol{\tau}_i)}}{\int e^{-\frac{1}{\lambda} \tilde{S}(\boldsymbol{\tau}_i)} d\boldsymbol{\tau}_i}, \quad \mathbf{u}_L(\boldsymbol{\tau}_i) = \frac{\mathbf{R}^{-1} \mathbf{g}_{t_i}^{(c)} \mathbf{g}_{t_i}^{(c)T}}{\mathbf{g}_{t_i}^{(c)T} \mathbf{R}^{-1} \mathbf{g}_{t_i}^{(c)}} \boldsymbol{\varepsilon}_{t_i},$$

# Iterative Path integrals with Parametrized policies

$$\mathbf{a}_{t_i} = \mathbf{g}_{t_i}^T (\boldsymbol{\theta} + \boldsymbol{\varepsilon}_{t_i})$$

$$\mathbf{u}_{t_i} = \int P(\boldsymbol{\tau}_i) \mathbf{u}_L(\boldsymbol{\tau}_i) d\boldsymbol{\tau}_i^{(c)},$$

$$P(\boldsymbol{\tau}_i) = \frac{e^{-\frac{1}{\lambda} \tilde{S}(\boldsymbol{\tau}_i)}}{\int e^{-\frac{1}{\lambda} \tilde{S}(\boldsymbol{\tau}_i)} d\boldsymbol{\tau}_i}, \quad \mathbf{u}_L(\boldsymbol{\tau}_i) = \frac{\mathbf{R}^{-1} \mathbf{g}_{t_i}^{(c)} \mathbf{g}_{t_i}^{(c)T}}{\mathbf{g}_{t_i}^{(c)T} \mathbf{R}^{-1} \mathbf{g}_{t_i}^{(c)}} \boldsymbol{\varepsilon}_{t_i}$$

$$\tilde{S}(\boldsymbol{\tau}_i) = \phi_{t_N} + \sum_{j=i}^{N-1} q_{t_j} + \frac{1}{2} \sum_{j=i}^{N-1} \boldsymbol{\varepsilon}_{t_j}^T \mathbf{M}_{t_j}^T \mathbf{R} \mathbf{M}_{t_j} \boldsymbol{\varepsilon}_{t_j}$$

$$\mathbf{M}_{t_j} = \frac{\mathbf{R}^{-1} \mathbf{g}_{t_j} \mathbf{g}_{t_j}^T}{\mathbf{g}_{t_j}^T \mathbf{R}^{-1} \mathbf{g}_{t_j}}$$

$$\begin{aligned} \boldsymbol{\theta}_{t_i}^{(new)} &= \int P(\boldsymbol{\tau}_i) \frac{\mathbf{R}^{-1} \mathbf{g}_{t_i} \mathbf{g}_{t_i}^T (\boldsymbol{\theta} + \boldsymbol{\varepsilon}_{t_i})}{\mathbf{g}_{t_i}^T \mathbf{R}^{-1} \mathbf{g}_{t_i}} d\boldsymbol{\tau}_i \\ &= \int P(\boldsymbol{\tau}_i) \frac{\mathbf{R}^{-1} \mathbf{g}_{t_i} \mathbf{g}_{t_i}^T \boldsymbol{\varepsilon}_{t_i}}{\mathbf{g}_{t_i}^T \mathbf{R}^{-1} \mathbf{g}_{t_i}} d\boldsymbol{\tau}_i + \frac{\mathbf{R}^{-1} \mathbf{g}_{t_i} \mathbf{g}_{t_i}^T \boldsymbol{\theta}}{\mathbf{g}_{t_i}^T \mathbf{R}^{-1} \mathbf{g}_{t_i}} \\ &= \delta \boldsymbol{\theta}_{t_i} + \frac{\mathbf{R}^{-1} \mathbf{g}_{t_i} \mathbf{g}_{t_i}^T}{\text{trace}(\mathbf{R}^{-1} \mathbf{g}_{t_i} \mathbf{g}_{t_i}^T)} \boldsymbol{\theta} \\ &= \delta \boldsymbol{\theta}_{t_i} + \mathbf{M}_{t_i} \boldsymbol{\theta}. \end{aligned}$$

Compare to DDP

# PI<sup>2</sup>

## Update step

$$P(\tau_i) = \frac{e^{-\frac{1}{\lambda}S(\tau_i)}}{\int e^{-\frac{1}{\lambda}S(\tau_i)} d\tau_i}, \quad (35)$$

$$S(\tau_i) = \phi_{t_N} + \sum_{j=i}^{N-1} q_{t_j} dt + \frac{1}{2} \sum_{j=i}^{N-1} (\theta + \mathbf{M}_{t_j} \boldsymbol{\varepsilon}_{t_j})^T \mathbf{R} (\theta + \mathbf{M}_{t_j} \boldsymbol{\varepsilon}_{t_j}) dt, \quad (36)$$

$$\delta\theta_{t_i} = \int P(\tau_i) \mathbf{M}_{t_i} \boldsymbol{\varepsilon}_{t_i} d\tau_i, \quad (37)$$

$$[\delta\theta]_j = \frac{\sum_{i=0}^{N-1} (N-i) w_{j,t_i} [\delta\theta_{t_i}]_j}{\sum_{i=0}^{N-1} w_{j,t_i} (N-i)}, \quad (38)$$

$$\theta^{(new)} = \theta^{(old)} + \delta\theta.$$

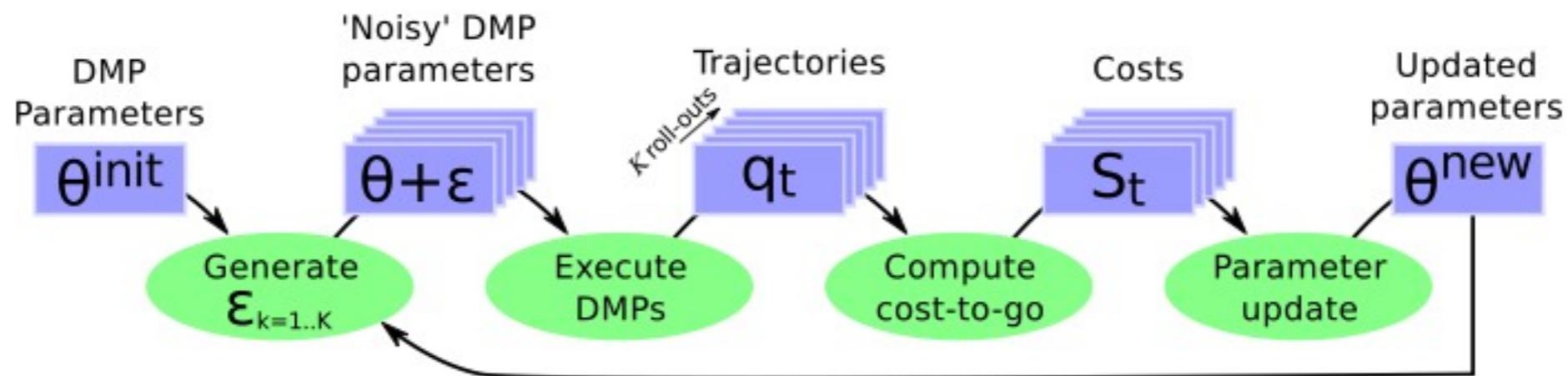


Fig. 2. Overview of the PI<sup>2</sup> algorithm.



# PI<sup>2</sup>

## Policy Improvement with Path Integrals

- Given:

- An immediate cost function  $r_t = q_t + \theta_t^T \mathbf{R} \theta_t$  (cf. 1)
- A terminal cost term  $\phi_{t_N}$  (cf. 1)
- A stochastic parameterized policy  $\mathbf{a}_t = \mathbf{g}_t^T (\theta + \epsilon_t)$  (cf. 25)
- The basis function  $\mathbf{g}_t$  from the system dynamics (cf. 3 and Section 2.5.1)
- The variance  $\Sigma_{\epsilon}$  of the mean-zero noise  $\epsilon_t$
- The initial parameter vector  $\theta$

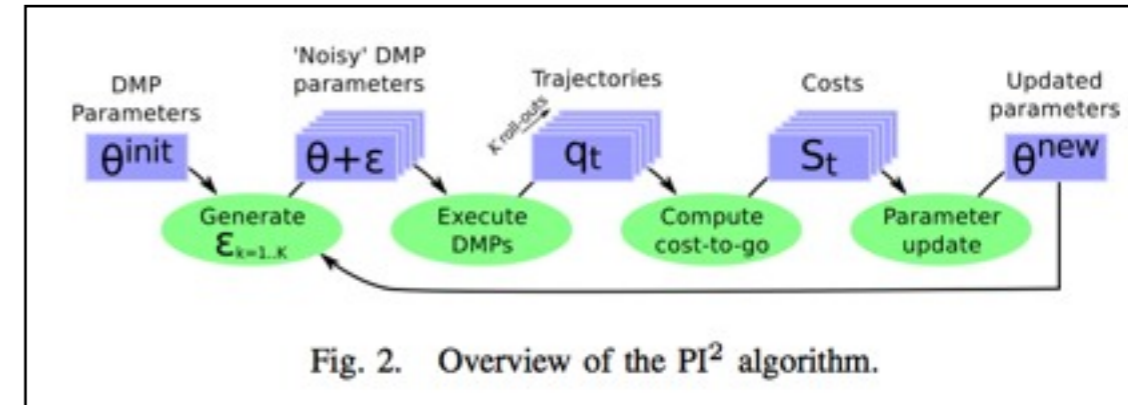


Fig. 2. Overview of the PI<sup>2</sup> algorithm.

- Repeat until convergence of the trajectory cost  $R$ :

- Create  $K$  roll-outs of the system from the same start state  $\mathbf{x}_0$  using stochastic parameters  $\theta + \epsilon_t$  at every time step
- For  $k = 1 \dots K$ , compute:
  - \*  $P(\tau_{i,k}) = \frac{e^{-\frac{1}{\lambda} S(\tau_{i,k})}}{\sum_{k=1}^K [e^{-\frac{1}{\lambda} S(\tau_{i,k})}]}$
  - \*  $S(\tau_{i,k}) = \phi_{t_N,k} + \sum_{j=i}^{N-1} q_{t_j,k} + \frac{1}{2} \sum_{j=i+1}^{N-1} (\theta + \mathbf{M}_{t_j,k} \epsilon_{t_j,k})^T \mathbf{R} (\theta + \mathbf{M}_{t_j,k} \epsilon_{t_j,k})$
  - \*  $\mathbf{M}_{t_j,k} = \frac{\mathbf{R}^{-1} \mathbf{g}_{t_j,k} \mathbf{g}_{t_j,k}^T}{\mathbf{g}_{t_j,k}^T \mathbf{R}^{-1} \mathbf{g}_{t_j,k}}$
- For  $i = 1 \dots (N-1)$ , compute:
  - \*  $\delta \theta_{t_i} = \sum_{k=1}^K [P(\tau_{i,k}) \mathbf{M}_{t_i,k} \epsilon_{t_i,k}]$
- Compute  $[\delta \theta]_j = \frac{\sum_{i=0}^{N-1} (N-i) w_{j,t_i} [\delta \theta_{t_i}]_j}{\sum_{i=0}^{N-1} w_{j,t_i} (N-i)}$
- Update  $\theta \leftarrow \theta + \delta \theta$
- Create one noiseless roll-out to check the trajectory cost  $R = \phi_{t_N} + \sum_{i=0}^{N-1} r_{t_i}$ . In case the noise cannot be turned off, that is, a stochastic system, multiple roll-outs need be averaged.

# Simplifications to PI2

$$S(\boldsymbol{\tau}_{i,k}) = \phi_{t_N,k} + \sum_{j=i}^{N-1} r_{t_j,k} +$$

~~$$\frac{1}{2} \sum_{j=i+1}^{N-1} (\boldsymbol{\theta} + \mathbf{M}_{t_j,k} \boldsymbol{\epsilon}^{\boldsymbol{\theta}}_{t_j,k})^T \mathbf{R} (\boldsymbol{\theta} + \mathbf{M}_{t_j,k} \boldsymbol{\epsilon}^{\boldsymbol{\theta}}_{t_j,k}) \quad (7)$$~~

~~$$\mathbf{M}_{t_j,k} = \frac{\mathbf{R}^{-1} \mathbf{g}_{t_j} \mathbf{g}_{t_j}^T}{\mathbf{g}_{t_j}^T \mathbf{R}^{-1} \mathbf{g}_{t_j}} \quad (8)$$~~

$$P(\boldsymbol{\tau}_{i,k}) = \frac{e^{-\frac{1}{\lambda} S(\boldsymbol{\tau}_{i,k})}}{\sum_{l=1}^K [e^{-\frac{1}{\lambda} S(\boldsymbol{\tau}_{i,l})}]} \quad (9)$$

~~$$\delta \boldsymbol{\theta}_{t_i} = \sum_{k=1}^K [P(\boldsymbol{\tau}_{i,k}) \mathbf{M}_{t_i,k} \boldsymbol{\epsilon}^{\boldsymbol{\theta}}_{t_i,k}] \quad (10)$$~~

~~$$[\delta \boldsymbol{\theta}]_j = \frac{\sum_{i=0}^{N-1} (N-i) w_{j,t_i} [\delta \boldsymbol{\theta}_{t_i}]_j}{\sum_{i=0}^{N-1} w_{j,t_i} (N-i)} \quad (11)$$~~

~~$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \delta \boldsymbol{\theta} \quad (12)$$~~



# Simplifications

$$M = I$$

Only use total cost to go

$$P(\tau_{0,k}) = \frac{e^{-\frac{1}{\lambda} J(\tau_{0,k})}}{\sum_{l=1}^K [e^{-\frac{1}{\lambda} J(\tau_{0,l})}]} \quad \text{Probability} \quad (21)$$

$$\delta g = \sum_{k=1}^K [P(\tau_{0,k}) \epsilon^g_k] \quad \text{Weighted averaging} \quad (22)$$

$$g \leftarrow g + \delta g \quad \text{Update} \quad (23)$$

# ... a few more things

Elitism: Remember overall best few and use  
in update

Lambda: Use schedule to 'freeze' the system

# Credits & Refs

Path Integral Based Stochastic Optimal Control for Rigid Body Dynamics  
E.A.Theodorou, J. Buchli and S. Schaal

A Generalized Path Integral Control Approach to Reinforcement Learning.  
Evangelos A.Theodorou, Jonas Buchli, Stefan Schaal

Learning Motion Primitive Goals for Robust Manipulation  
Freek Stulp, Evangelos Theodorou, Mrinal Kalakrishnan, Peter Pastor, Ludovic Righetti, Stefan Schaal

Feynman Lectures on Physics



# EOF L7

