

# AA228: Policy Search

---

So far, we've seen how exact solution methods can be used to solve for a policy in an offline method. We've also seen how online planning can handle large state spaces by reasoning over actions from an initial state. In this notebook, we'll discuss **policy search**, which involves searching over the space of policy parameters rather than actions. This idea will carry forward through the next few lectures.

You can find this notebook in the github repo [here](#)

## Table of Contents

### AA228: Policy Search

- Motivation
- MDP Formulation: Inverted Pendulum
- Policy Parameterization
- Policy Evaluation
- Policy Search Overview
- Local Search (Hooke-Jeeves)
- Genetic Algorithms
- Cross Entropy Method
- Algorithm Comparison
- Next Steps: Gradient Information

## Motivation

---

**Q: Given the solution methods we've seen so far, how can we solve an MDP with continuous states or actions?**

- One option is discretization, however we risk losing fidelity
- In value function based methods, we can use approximate value iteration.
- A common approach is to parameterize a policy. Then we can optimize over the policy parameters, which are usually much lower dimension than the action space.

Policy search helps us scale to very large or continuous state spaces.

```
1 begin
2     using POMDPs
3     using POMDPTools
4     using Distributions
5     using Parameters
6     using Random
7 end
```

## MDP Formulation: Inverted Pendulum

Let's define a running example with continuous states and actions to get started.

The inverted pendulum problem involves stabilizing a pendulum in an inverted position. Suppose we have a motor mounted at the pendulum's pivot point. The objective is to control the angle of the pendulum so that it stays as close to vertical as possible despite any disturbances that may occur, such as gusts of wind. This problem is challenging because the inverted pendulum is inherently unstable, making it difficult to maintain its balance.

Let's define the angle of the pendulum from vertical as  $\phi$ .

We'll assume the pendulum has some fixed length  $l$  and mass  $m$ .

**Q: What should we use as the MDP state? Actions?**

A:

- State: Angle information  $\phi$  and  $\dot{\phi}$
- Action: Motor torque

```
1 md"""
2 A:
3
4 * State: Angle information  $\phi$  and  $\dot{\phi}$ 
5 * Action: Motor torque
6 """
```

**Q: What is the transition model? (The actual equations are not important, but describe in words.)**

A:

- Physics based model

```
1 md"""
2 A:
3 * Physics based model
4 """
```

**Q: What components might our reward function have?**

A:

- Penalize angles far from vertical
- Penalize based on angular velocity
- Penalize large torques

```
1 md"""
2 A:
3 * Penalize angles far from vertical
4 * Penalize based on angular velocity
5 * Penalize large torques
6 """
```

Great! Now let's look at how we can set define this MDP in code. We'll use the POMDPs.jl environment.

First, we'll define a struct type. This is a container for useful data about the pendulum.

PendulumMDP

```
1 @with_kw struct PendulumMDP <: MDP{Array{Float64}, Array{Float64}}
2     Rstep = 1 # Reward earned on each step of the simulation
3     λcost = 1 # Coefficient to the traditional OpenAIGym Reward
4     max_speed::Float64 = 8.
5     max_torque::Float64 = 100.
6     dt::Float64 = .05
7     g::Float64 = 10.
8     m::Float64 = 1.
9     l::Float64 = 1.
10    γ::Float64 = 0.99
11 end
```

```
1 POMDPs.discount(mdp::PendulumMDP) = mdp.γ
```

For many problems, explicitly writing the transition model  $T(s' | s, a)$  and reward function  $R(s, a)$  can be difficult. Here, we define a **generative model** of the dynamics and reward.

```

1 function pendulum_dynamics(env, s, a; rng::AbstractRNG = Random.GLOBAL_RNG)
2     θ, ω = s[1], s[2]
3     dt, g, m, l = env.dt, env.g, env.m, env.l
4
5     a = a[1]
6     a = clamp(a, -env.max_torque, env.max_torque)
7     costs = angle_normalize(θ)^2 + 0.1f0 * ω^2 + 0.001f0 * a^2
8
9     ω = ω + (-3. * g / (2 * l) * sin(θ + π) + 3. * a / (m * l^2)) * dt
10    θ = angle_normalize(θ + ω * dt)
11    ω = clamp(ω, -env.max_speed, env.max_speed)
12
13    sp = [θ, ω]
14    r = env.Rstep - env.λcost*costs
15    return sp, r
16 end;

```

```

1 angle_normalize(x) = mod((x+π), (2*π)) - π;

```

```

1 function POMDPs.gen(mdp::PendulumMDP, s, a, rng::AbstractRNG = Random.GLOBAL_RNG)
2     sp, r = pendulum_dynamics(mdp, s, a, rng=rng)
3     (sp = sp, r = r)
4 end

```

Define an initial state distribution

```

1 function POMDPs.initialstate(mdp::PendulumMDP)
2     θ0 = Distributions.Uniform(-π/6., π/6.)
3     ω0 = Distributions.Uniform(-0.1, 0.1)
4     ImplicitDistribution((rng) -> [rand(rng, θ0), rand(rng, ω0)])
5 end

```

Rendering functions are below (hidden cells). Feel free to look under the hood if you are curious!

# Policy Parameterization

Now that we've established the components of our MDP, let's start to think about how to solve it. We already discussed how offline methods and online tree-search methods might be difficult to apply.

We introduce the notion of a **parameterized policy**. We can denote the action of policy  $\pi$  at state  $s$  parameterized by  $\theta$  as

$$a = \pi_{\theta}(s)$$

for deterministic policies, and

$$a \sim \pi_{\theta}(a \mid s)$$

for stochastic policies.

**Policy space is often lower-dimensional than state space, and can often be searched more easily.**

The parameters  $\theta$  may be a vector or some other more complex representation. For example, we may want to represent our policy using a neural network with a particular structure. We would use  $\theta$  to represent the weights in the network.

**Q: How could we parameterize a policy for the inverted pendulum problem? Assume our state vector is  $s = [\phi, \dot{\phi}]$**

**A:** One option is a weighted combination of elements in the state vector

$$\pi_{\theta}(s) = \theta_1 \phi + \theta_2 \dot{\phi} = \theta^T s$$

This is a deterministic policy. If we wanted a stochastic policy, we could also use a linear gaussian model.

$$\pi_{\theta}(a \mid s) = \mathcal{N}(a \mid \theta_A s + \theta_b, \theta_{\Sigma})$$

# Policy Evaluation

The expected discounted return of a policy  $\pi$  from initial state distribution  $b(s)$  is

$$U(\pi) = \sum_s b(s) U^\pi(s)$$

When we have a large or continuous state space, we often cannot compute the utility of following a policy  $U(\pi)$  exactly. Instead, we rewrite  $U(\pi)$  in terms of trajectories of states, actions, and rewards under the policy  $\pi$ .

The key idea is that we want to estimate utility based on *simulated trajectories*,  $\tau$ . We call the sum of rewards for trajectory  $\tau$ ,  $R(\tau)$  the *trajectory return*. Now, we can write the utility of following policy  $\pi$  as:

$$U(\pi) = \mathbb{E}[R(\tau)]$$

The expected value (or mean) return can be *approximated* by taking the mean total reward of many trajectories:

$$U(\pi) \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$$

This is sometimes called Monte Carlo policy evaluation.

The POMDPs.jl package provides us with a convenient way to get the sum of discounted returns from a simulation (or a 'rollout'). The **RolloutSimulator** type simulates a given policy for a fixed number of steps and returns the sum of discounted returns. Once we create a simulator like:

```
sim = RolloutSimulator(max_steps=max_steps)
```

We can compute get the total discounted reward using **R=simulate(sim, mdp, policy)**

**Let's write a function to perform Monte Carlo policy evaluation.**

```
1 function mc_policy_evaluation(mdp::MDP, pi::Policy; m=100, max_steps=100)
2     sim = RolloutSimulator(max_steps=max_steps)
3     return mean([simulate(sim, mdp, pi) for _=1:m])
4 end;
```

# Policy Search Overview

In policy search, our goal is to optimize a policy's utility with respect to its parameters. In other words, we search over the parameter space for a set of parameters that maximize our utility.

Here, we'll try out policy search for a simple 2D policy parameterization:

$$\pi_{\theta}(s) = \theta_1 s_1 + \theta_2 s_2 = \theta^T s$$

First, let's create our MDP

```
mdp = PendulumMDP
  Rstep: Int64 1
  λcost: Int64 1
  max_speed: Float64 8.0
  max_torque: Float64 100.0
  dt: Float64 0.05
  g: Float64 10.0
  m: Float64 1.0
  l: Float64 1.0
  γ: Float64 0.99
```

```
1 mdp = PendulumMDP()
```

Select  $\theta_1$  and  $\theta_2$  to maximize the utility

theta1 :  theta2 : 

```
θ = [-30.0, -30.0]
```

```
1 θ = [θ1, θ2]
```

```
1 policy = FunctionPolicy((s) -> [θ' * s]);
```

```
mean_utility = -1165.8407861514054
```

```
1 mean_utility = mc_policy_evaluation(mdp, policy)
```

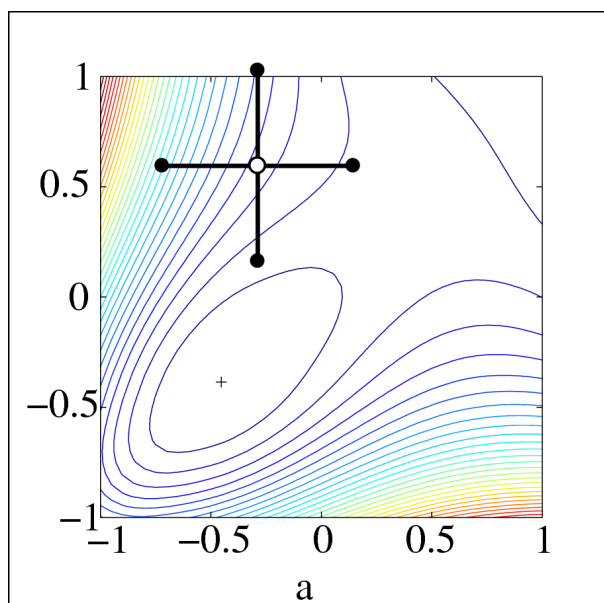
```
1 # begin
2 #   animate_pendulum(mdp, policy, "toyexample.gif")
3 #   LocalResource("toyexample.gif")
4 # end
```

# Local Search (Hooke-Jeeves)

Hooke-Jeeves is an example of a local search method. We saw a local search algorithm in structure learning. In local search, we evaluate the neighbors of a current design point ( $\theta$ ), and move to the neighbor that most improves the value.

The Hooke-Jeeves algorithm is one example of local search. The algorithm takes a step of size  $\pm\alpha$  in each of the *coordinate directions* from the current  $\theta$ . If an improvement is found, it moves to the best point. If no improvement is found, the algorithm shrinks the step size  $\alpha$  and repeats. The algorithm terminates once  $\alpha$  reaches some minimum value.

The algorithm is illustrated below



Now let's code up Hooke-Jeeves and apply it to the inverted pendulum! This code very closely follows the textbook. First, we'll define a special struct to store useful attributes of the algorithm.

```
1 struct HookeJeevesPolicySearch
2      $\theta$  # initial parameterization
3      $\alpha$  # step size
4     c # step size reduction factor
5      $\epsilon$  # termination step size
6 end
```

Next, let's define the optimization procedure. We'll take in the algorithm attributes, as well as a function that takes in a parameter vector  $\theta$  and returns a Monte Carlo estimate of the expected utility.



```

1 function optimize(M::HookeJeevesPolicySearch, U)
2      $\theta$ ,  $\theta'$ ,  $\alpha$ , c,  $\epsilon$  = copy(M. $\theta$ ), similar(M. $\theta$ ), M. $\alpha$ , M.c, M. $\epsilon$ 
3     u = U( $\theta$ )
4     n = length( $\theta$ )
5     history = [copy( $\theta$ )]
6     while  $\alpha$  >  $\epsilon$ 
7         copyto!( $\theta'$ ,  $\theta$ )
8         best = (i=0, sgn=0, u=u)
9         for i in 1:n
10             for sgn in (-1,1)
11                  $\theta'$ [i] =  $\theta$ [i] + sgn* $\alpha$ 
12                 u' = U( $\theta'$ )
13                 if u' > best.u
14                     best = (i=i, sgn=sgn, u=u')
15                 end
16             end
17              $\theta'$ [i] =  $\theta$ [i]
18         end
19         if best.i != 0
20              $\theta$ [best.i] += best.sgn* $\alpha$ 
21             u = best.u
22         else
23              $\alpha$  *= c
24         end
25         push!(history, copy( $\theta$ ))
26     end
27     return  $\theta$ , history
28 end;

```

Now we'll define our function  $U(\theta)$

```

1 U( $\theta$ ) = mc_policy_evaluation(mdp, FunctionPolicy((s)->[ $\theta$ '*s]), max_steps=500, m=10);

```

Specify the attributes of our algorithm

```

1 hookejeeves = HookeJeevesPolicySearch(rand(2), 1.0, 0.5, 1e-1);

```

Now run the optimization!

```

1  $\theta_{hj}$ , history = optimize(hookejeeves, U);

```

```
[-5.74906, -0.282411]
```

```
1  $\theta_{hj}$ 
```

```
97.3306550038216
```

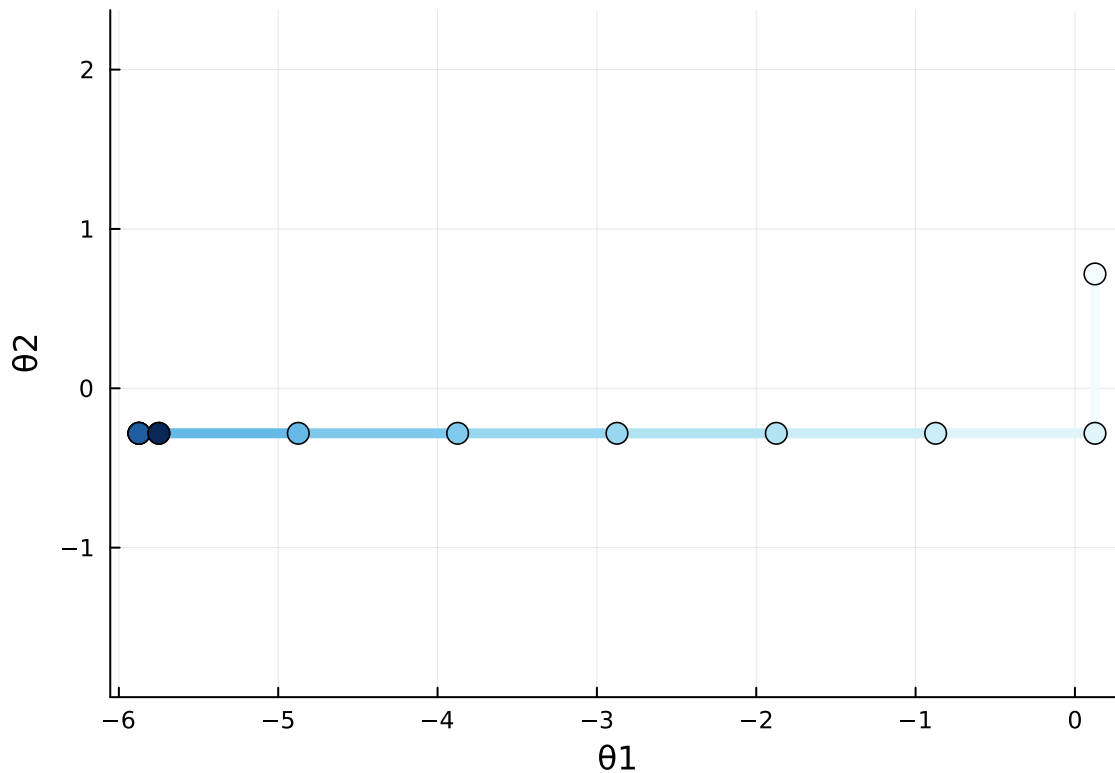
```
1 U( $\theta_{hj}$ )
```

```

1 # begin
2 #    $\pi_{cem} = \text{FunctionPolicy}((s) \rightarrow [\theta_{hj}' * s])$ 
3 #   animate_pendulum(mdp,  $\pi_{cem}$ , "hj.gif")
4 #   LocalResource("hj.gif")
5 # end

```

We can also examine the history of points evaluated by the algorithm. The algorithm starts at light blue and progresses toward darker blue.



**Q:** Suppose we are performing Hooke-Jeeves for policy parameterized by a 3-element vector. We are currently evaluating the point  $[2, -5, 10]$ . Assume the step size is 1.

**What points will Hooke-Jeeves evaluate next?**

**A:** Hooke-Jeeves will evaluate **6** points

- Coordinate 1:  $[3, -5, 10]$  and  $[1, -5, 10]$
- Coordinate 2:  $[2, -4, 10]$  and  $[2, -6, 10]$
- Coordinate 3:  $[2, -5, 11]$  and  $[2, -5, 9]$

**Q:**Hooke-Jeeves evaluates each of these points, and finds they have utilities

(10, 23, 16, 5, 34, 27)

**. The current policy has a value of 35. What will Hooke-Jeeves do?**

A: Hooke-Jeeves will contract the step size. The utility at each neighbor is lower than the current policy.

## Genetic Algorithms

Local search algorithms can easily become stuck in local minima. Population-based algorithms maintain a set of points in the parameter space. By maintaining and changing the set of parameters, population-based algorithms can be less susceptible to becoming stuck. However, they are not guaranteed to converge to the global optimum.

Genetic algorithms are population-based algorithms that are inspired by biological evolution. The algorithm starts with a population of points  $m$  in parameter space, called 'individuals':  $\theta^{(1)}, \dots, \theta^{(m)}$ . We compute  $U(\theta)$  for each of point. The top-performing samples, are called *elite samples*. At the next iteration, some number  $m_{elite}$  of the elite samples are chosen. New points in the population are created by adding gaussian noise to the elite individuals.

```
1 struct GeneticPolicySearch
2      $\theta$ s # initial population
3      $\sigma$  # initial standard deviation
4     m_elite # number of elite samples
5     k_max # number of iterations
6 end
```

```
1 function optimize(M::GeneticPolicySearch, U)
2      $\theta$ s,  $\sigma$  = M. $\theta$ s, M. $\sigma$ 
3     n, m = length(first( $\theta$ s)), length( $\theta$ s)
4     history = []
5     for k in 1:M.k_max
6         us = [U( $\theta$ ) for  $\theta$  in  $\theta$ s]
7         sp = sortperm(us, rev=true)
8          $\theta\_best$  =  $\theta$ s[sp[1]]
9         push!(history, (copy( $\theta$ s), copy( $\theta$ s[sp[1:M.m_elite]])))
10        rand_elite() =  $\theta$ s[sp[rand(1:M.m_elite)]]
11         $\theta$ s = [rand_elite() +  $\sigma$ .*randn(n) for i in 1:(m-1)]
12        push!( $\theta$ s,  $\theta\_best$ )
13    end
14    return last( $\theta$ s), history
15 end;
```

## Creating an initial population

$[[16.5313, 21.5613], [2.23064, -10.6746], [22.4038, -18.8571], [23.337, -5.95289], [3.8042$

```
1 begin
2     npop = 50
3      $\theta\theta$  = [50 .* rand(2).-25 for _=1:npop]
4 end
```

```
1 ga = GeneticPolicySearch( $\theta\theta$ , 1.0, 10, 10);
```

```
1  $\theta$ ga, hga = optimize(ga, U);
```

The final parameters found using the genetic algorithm are:

$[-9.62394, -2.95826]$

```
1  $\theta$ ga
```

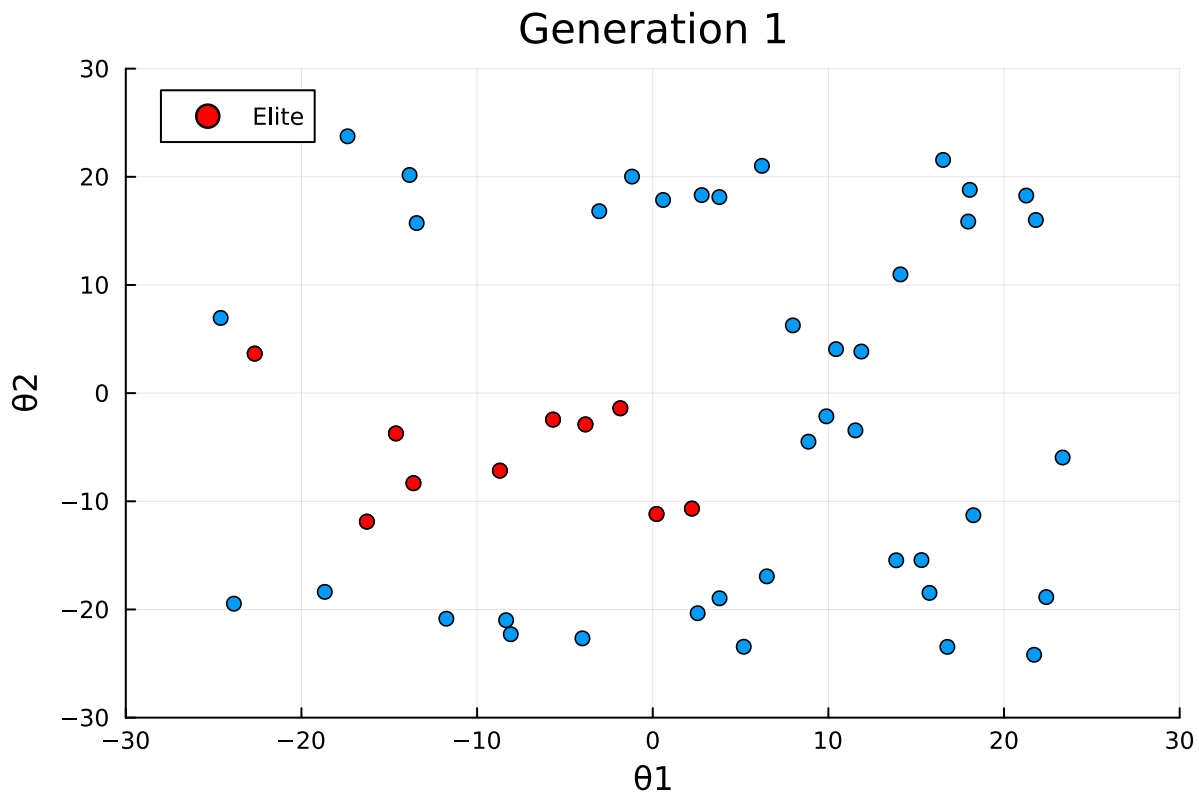
With an expected utility of:

98.68047010966082

```
1 U( $\theta$ ga)
```

Let's take a look at the population over each iteration. Try playing with the initial population, standard deviation, and other parameters and see how the algorithm behaves.





**Q: How will the solution found by a genetic algorithm depend on the initial population? How would it depend on the magnitude of noise added to elite samples?**

**A:** Generally, a more distributed population can more effectively explore the parameter space at the cost of computation (need for policy evaluation for more individuals) and convergence speed.

There is a similar story for the magnitude for the random perturbations. Adding random noise helps explore the policy space, but can also slow down convergence.

**Q: Is a genetic algorithm guaranteed to converge to the optimal policy?**

# Cross Entropy Method

The cross entropy method involves updating a search distribution over the parameters. The distribution over parameters  $p(\theta \mid \psi)$  has its own parameters  $\psi$ . Typically, we use a Gaussian distribution, where  $\psi$  represents the mean and covariance matrix.

The algorithm iteratively updates the parameters  $\psi$ . At each iteration, we draw  $m$  samples from the associated distribution and then update  $\psi$  to fit a set of elite samples. We stop after a fixed number of iterations, or when the search distribution becomes very focused.

```
1 struct CrossEntropyPolicySearch
2     p # initial distribution
3     m # number of samples
4     m_elite # number of elite samples
5     k_max # number of iterations
6 end
```

```
1 function optimize_dist(M::CrossEntropyPolicySearch, U)
2     p, m, m_elite, k_max = M.p, M.m, M.m_elite, M.k_max
3     history = []
4     for k in 1:k_max
5         θs = rand(p, m)
6         us = [U(θs[:,i]) for i in 1:m]
7         θ_elite = θs[:,sortperm(us)[(m-m_elite+1):m]]
8         push!(history, (p, copy(θs), copy(θ_elite)))
9         p = Distributions.fit(typeof(p), θ_elite)
10
11     end
12     return p, history
13 end;
```

```
1 function optimize(M, U)
2     d, history = optimize_dist(M, U)
3     return Distributions.mode(d), history
4 end;
```

A key step of using the algorithm is selecting the *initial distribution*. The distribution should cover the parameter space of interest.

```
p0 = DiagonalNormal(
    dim: 2
    μ: [0.0, 0.0]
    Σ: [100.0 0.0; 0.0 100.0]
)
```

```
1 p0 = MvNormal(zeros(2), [10, 10])
```

```
cem = CrossEntropyPolicySearch(DiagonalNormal(
    dim: 2
    μ: [0.0, 0.0]
    Σ: [100.0 0.0; 0.0 100.0]
), 50, 10, 10)
```

```
1 cem = CrossEntropyPolicySearch(p0, 50, 10, 10)
```

```
1  $\theta_{cem}$ ,  $h_{cem}$  = optimize(cem, U);
```

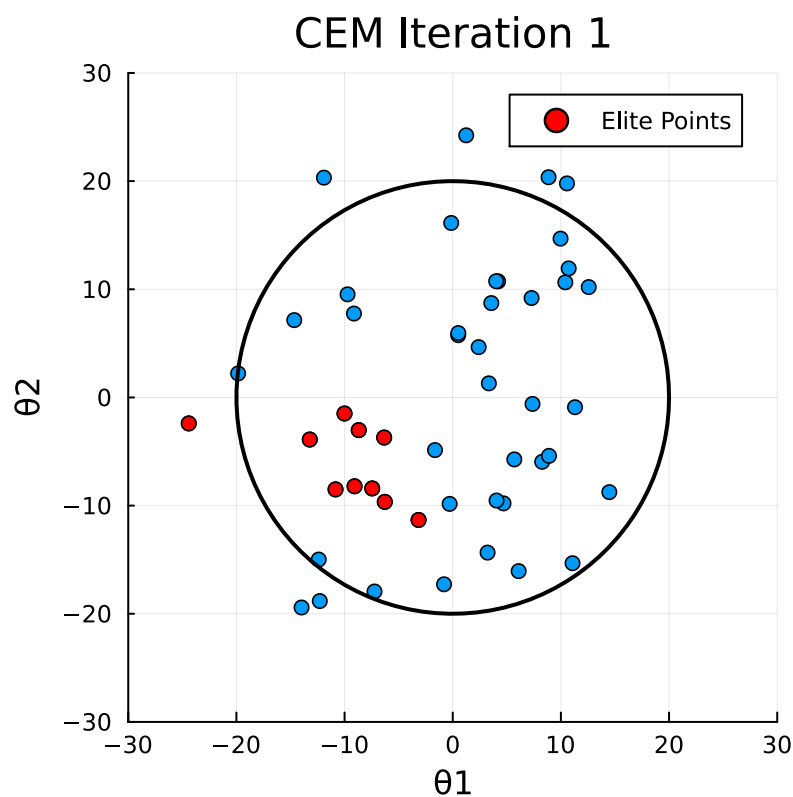
```
[-15.8694, -5.4521]
```

```
1  $\theta_{cem}$ 
```

```
98.47451969676476
```

```
1  $U(\theta_{cem})$ 
```

We can also examine how the proposal distribution changes over iterations.



**Q:** We are performing the Cross Entropy Method on a 1D parameter space. We have elite samples at (1, 4, 2.5, 10) What will be the updated parameters for our Gaussian search distribution?

Recall the maximum likelihood estimate for Gaussian's parameters with samples  $o_1, \dots, o_m$ :

$$\hat{\mu} = \frac{\sum_i o_i}{m}$$

$$\hat{\sigma}^2 = \frac{\sum_i (o_i - \hat{\mu})^2}{m}$$

A: We use a maximum likelihood estimate to update the parameters. For a Gaussian distribution and samples  $\theta_1$

The new **mean** is  $(1+4+2.5+10)/4 = 4.375$

The new **variance** is  $((1-4.375)^2 + (4-4.375)^2 + (2.5-4.375)^2 + (10-4.375)^2)/4 = 11.671875$

## Algorithm Comparison

Let's compare the performance of each algorithm on the inverted pendulum

```
1 begin
2   @show U( $\theta_j$ )
3   @show U( $\theta_{ga}$ )
4   @show U( $\theta_{cem}$ )
5 end;
```

```
U( $\theta_j$ ) = 98.40847554833884
U( $\theta_{ga}$ ) = 98.95232447297779
U( $\theta_{cem}$ ) = 98.75504052014878
```



They all do pretty well!

**Q:** What if our policy had many more parameters, say 100. Which algorithm would you pick?

A:

- A genetic algorithm might be a good choice here, although it would also require a large population to explore the parameter space. The cross entropy method is also a good choice, in higher dimensions more samples are required to fit the search distribution.
- Hooke Jeeves could be an OK option, but is subject to getting stuck in local optima.



**Q: Suppose that we want to perform policy optimization on a problem where we know that policies far apart in parameter space can have similar high utility. What are the advantages of genetic algorithms over hooke-jeeves? What about the Cross Entropy method?**

A: Here we would prefer population-based algorithms or the cross-entropy method. We can argue that both of these algorithms could be likely to find a global optimum. Genetic algorithms can effectively explore multimodal objectives. The cross entropy method can also represent multimodal objectives, and is less likely to converge to a local optima than Hooke-Jeeves

```
1 md"""
2 A:
3 Here we would prefer population-based algorithms or the cross-entropy method. We can
  argue that both of these algorithms could be likely to find a global optimum. Genetic
  algorithms can effectively explore multimodal objectives. The cross entropy method
  can also represent multimodal objectives, and is less likely to converge to a local
  optima than Hooke-Jeeves
4 """
```

# Next Steps: Gradient Information

Thus far, all of the algorithms for policy search have not used gradient information of the expected utility with respect to policy parameters. It turns out that optimizing policies with many parameters can be done much more efficiently with gradient information

How can we compute the gradient  $\nabla U(\theta)$ ?

One option is to use *finite differences*. The idea of finite differences comes from the linear approximation of the gradient. We can estimate the gradient (or the slope) of a function in 1D by checking how much the value of  $f(x)$  changes for some small change in  $x$ .

$$\frac{df}{dx} \approx \frac{f(x + \delta) - f(x)}{\delta}$$

If we extend this same idea to our utility function,

$$\nabla U(\theta) \approx \left[ \frac{U(\theta + \delta \mathbf{e}^1) - U(\theta)}{\delta}, \dots, \frac{U(\theta + \delta \mathbf{e}^n) - U(\theta)}{\delta} \right]$$

Where  $\mathbf{e}^i$  is the standard basis, which is zero everywhere except the  $i$ th component.

Luckily, there are some great packages in most programming languages that do this for us!

However, recall that our Monte Carlo estimate of the expected utility is stochastic, or noisy. This means that gradients of the utility function will have noise too! If the gradients are too noisy, they will provide very poor guidance for our policy search.

A key challenge in policy gradient estimation is dealing with noisy policy gradients.

```
1 using FiniteDiff
```

```
1 function deterministic_policy_evaluation(mdp::MDP, π::Policy; m=100, max_steps=500)
2     sim = RolloutSimulator(rng=MersenneTwister(42), max_steps=max_steps)
3     return mean([simulate(sim, mdp, π) for _=1:m])
4 end;
```

```
1 Ufd(θ) = deterministic_policy_evaluation(mdp, FunctionPolicy((s)->[θ'*s])),
    max_steps=500, m=5);
```

Let's try taking the policy gradient.

```
[310.228, -1280.2]
```

```
1 FiniteDiff.finite_difference_gradient(Ufd, [0.1, 0.1])
```

Now that we have an estimate of the gradient, how can we use it to improve the policy?

One of the simplest approaches is **gradient ascent**. Gradient ascent takes steps in parameter space along the gradient direction. The step size  $\alpha$  determines how far along the gradient direction the update moves. The update for  $\theta$  is

$$\theta \leftarrow \theta + \alpha \nabla U(\theta)$$

Determining the step size is a major challenge. Large steps can lead to faster progress to the optimum, but they can overshoot.

Let's try running a very simple version of gradient descent and see how it performs!

```
1 begin
2    $\theta_i$  = rand(2)
3    $\alpha$  = 1e-2
4   niter = 100
5   for k=1:niter
6     gradU = FiniteDiff.finite_difference_gradient(Ufd,  $\theta_i$ )
7      $\theta_i$  +=  $\alpha$  .* gradU
8   end
9 end
```

The final parameters are a little different than what was previously found. Why could that be?

```
[-6.21113, -0.739369]
```

```
1  $\theta_i$ 
```

```
97.60565017114206
```

```
1 Ufd( $\theta_i$ )
```

```
1 # begin
2 #    $\pi_{sgd}$  = FunctionPolicy((s)->[ $\theta_i$ '*s])
3 #   animate_pendulum(mdp,  $\pi_{sgd}$ , "sgd.gif")
4 #   LocalResource("sgd.gif")
5 # end
```

Gradient descent is able to find a decent policy pretty quickly! However, it isn't as good as the policies we've found previously. It turns out there are **much** better ways to estimate the gradient of a policy, and much more intelligent variations on gradient ascent. We'll explore these in the coming lecture.

This wraps up our discussion on policy search. I hope it was helpful!

