



RL Part 3.2: Probabilistic Policy Search

Gerhard Neumann
Jan Peters



What we have seen from the policy gradients

- Policy Search is a powerful and practical alternative to value function and model-based methods.
- Policy gradients have dominated this area for a long time and solidly working methods exist.
- Say still need a lot of samples and **we need to tune the learning rate**
- Learning the **exploration rate** is still an open problem



Outline of the Lecture

1. Introduction
2. Policy Updates by Weighted Maximum Likelihood
3. Relative Entropy Policy Search (REPS)
4. REPS for Contextual Policy Search
5. Conclusion

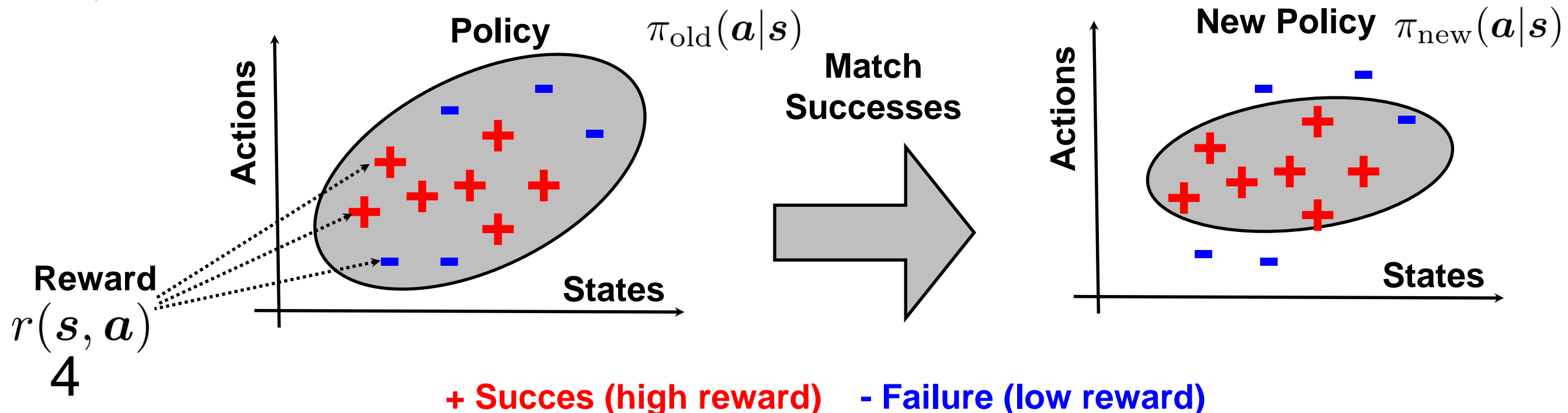


Success Matching Principle

“When learning from a set of their own trials in iterated decision problems, humans attempt to match **not the best taken action** but the **reward-weighted frequency** of their actions and outcomes” (Arrow, 1958).

- Why? We still need to explore!
- Create policies such that $\pi_{\text{new}}(\mathbf{a}|\mathbf{s}) \propto \pi_{\text{old}}(\mathbf{a}|\mathbf{s})r(\mathbf{s}, \mathbf{a})$

➡ Only possible for non-negative reward functions $r(\mathbf{s}, \mathbf{a})$



Reward
 $r(\mathbf{s}, \mathbf{a})$
4

Quick Recap: Episode-based Policy Search



For now, we will consider the **episode-based setting**

The policy evaluation strategy is needed to assess the quality of the samples

Episode-Based:

We directly assess the quality of a parameter vector $\theta^{[i]}$

$$R_{[i]} = \sum_{t=1}^T r_t^{[i]}$$

Data-set used for policy update

$$\mathcal{D}_{\text{episode}} = \left\{ \theta^{[i]}, R^{[i]} \right\}_{i=1 \dots N}$$

One data-point per trajectory

Episode-based policy evaluation strategy



The policy evaluation strategy is needed to assess the quality of the samples

Upper-level Policy :

We typically learn a distribution $\pi(\theta; \omega)$ over the parameters of low-level control policy $\pi(a|s; \theta)$

$\pi(\theta; \omega)$ is called **upper-level policy, e.g.** $\mathcal{N}(\theta|\mu, \Sigma)$

ω ... parameters of upper level policy

To reduce variance in the returns, $\pi(a|s; \theta)$ is often modelled as deterministic policy, i.e.,

$$\pi(a|s; \theta) \rightarrow a = \pi(s)$$

Works for a **moderate number of parameters (e.g. DMPs)**



Policy Update by Success Matching

Iterate:

Sample state-actions with current policy $\theta^{[i]} \sim \pi(\theta; \omega_k)$

Compute Target Distribution: „Success“ weighted policy on the samples

$$\tilde{\pi}(\theta^{[i]}) \propto w^{[i]} \pi_k(\theta^{[i]}; \omega_k) \quad w^{[i]} = f(R^{[i]})$$

We need to transform the reward with f in **a non-negative weight** (improper probability distribution)

Fit new parametric policy $\pi(\theta^{[i]}; \omega_{k+1})$ to target distribution

$$\omega_{k+1} = \operatorname{argmax}_{\omega} \sum_i w^{[i]} \log \pi(\theta^{[i]}; \omega)$$

➡ **Weighted Maximum Likelihood Estimate**

Policy Updates by Weighted ML

$$\omega_{k+1} = \operatorname{argmax}_{\omega} \sum_i w^{[i]} \pi(\theta^{[i]}; \omega)$$

Why is it cool?

No learning rate involved

Can be **computed efficiently** for many distributions

We can directly „jump“ to the desired distribution

Why can we do a **weighted ML estimate** with $w^{[i]}$ as weights?

For now, we will assume that the weights $w^{[i]}$ are given



Outline of the Lecture

1. Introduction
- 2. Policy Updates by Weighted Maximum Likelihood**
3. Relative Entropy Policy Search (REPS)
4. REPS for Contextual Policy Search
5. Conclusion

Policy Updates by Weighted ML

Why can we do a **weighted ML estimate** with $w^{[i]}$ as weights?

Problem: We want to find a parametric distribution $\pi(\boldsymbol{\theta}; \boldsymbol{\omega})$ that best fits the distribution $\tilde{\pi}(\boldsymbol{\theta}^{[i]}) \propto w^{[i]} \pi(\boldsymbol{\theta}^{[i]}; \boldsymbol{\omega}_k)$,

We can do that by minimizing the expected KL between $\tilde{\pi}(\boldsymbol{\theta}^{[i]})$

and $\pi(\boldsymbol{\theta}; \boldsymbol{\omega})$

**Minimizing the KL
is the same as maximizing
the likelihood!!**

$$\boldsymbol{\omega}_{k+1} = \operatorname{argmin}_{\boldsymbol{\omega}} \operatorname{KL}(\tilde{\pi}(\boldsymbol{\theta}^{[i]}) || \pi(\boldsymbol{\theta}^{[i]}; \boldsymbol{\omega}))$$

$$= \operatorname{argmin}_{\boldsymbol{\omega}} \int \tilde{\pi}(\boldsymbol{\theta}) \log \frac{\tilde{\pi}(\boldsymbol{\theta})}{\pi(\boldsymbol{\theta}; \boldsymbol{\omega})} d\boldsymbol{\theta}$$

$$\approx \operatorname{argmax}_{\boldsymbol{\omega}} \sum_i \frac{\tilde{\pi}(\boldsymbol{\theta}^{[i]})}{\pi(\boldsymbol{\theta}^{[i]}; \boldsymbol{\omega}_k)} \log \pi(\boldsymbol{\theta}^{[i]}; \boldsymbol{\omega})$$

$$\boldsymbol{\omega}_{k+1} = \operatorname{argmax}_{\boldsymbol{\omega}} \sum_i w^{[i]} \log \pi(\boldsymbol{\theta}^{[i]}; \boldsymbol{\omega})$$

We sampled from
the old policy

Weighted Maximum Likelihood Solutions...



If the upper-level policy $\pi(\boldsymbol{\theta}; \boldsymbol{\omega})$ is **Gaussian** $\mathcal{N}(\boldsymbol{\theta} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$ then mean and covariance are given by:

$$\boldsymbol{\mu} = \frac{\sum_i w^{[i]} \mathbf{a}^{[i]}}{\sum_i w^{[i]}}$$

Weighted mean

$$\boldsymbol{\Sigma} = \frac{\sum_i w^{[i]} (\mathbf{a}^{[i]} - \boldsymbol{\mu})^T (\mathbf{a}^{[i]} - \boldsymbol{\mu})}{\sum_i w^{[i]}}$$

Weighted covariance

Full covariance matrix  **correlated exploration** in parameter space

But more general: Also for mixture models, GPs and so on...

Weighted Maximum Likelihood Solutions...



So **where are the weights** $w^{[i]} = f(R^{[i]})$ coming from?

We need to transform the returns in an **improper probability distribution**

Expectation-Maximization Based Algorithms

One way of deriving the weighted ML updates

EM algorithms introduce a reward event C

$$w^{[i]} = p(C = 1 | \tau^{[i]}) \propto \exp(\beta R^{[i]})$$

Hence, the weight is given by an **exponential transformation of the return**

Expectation Maximization



Some notes on Expectation-Maximization in this context

EM is a method for Max. Likelihood in the case of latent (unobserved) variables

Observed variable: Reward Event $C = 1$ (we want to get reward)

Unobserved variable: Trajectory τ (or parameters) that created the reward event

E-Step: Estimate new desired distribution

M-Step: Estimate new policy parameters from the weighted samples

Step-based Policy Search Versions of EM: PoWER (Kober 2008), Reward-Weighted Regression (Peters 2007)

Expectation Maximization



Some notes on Expectation-Maximization in this context

Formally, the reward transformation is hard to motivate

$$w^{[i]} = p(C = 1 | \tau^{[i]}) \propto \exp(\beta R^{[i]})$$

β ... temperature parameter. Needs to be hand-tuned (task specific)

In stochastic environments, we do not optimize the expected reward any more as...

$$\mathbb{E}_{p(\tau)} [\exp(R(\tau))] \neq \exp(\mathbb{E}_{p(\tau)} [R(\tau)])$$

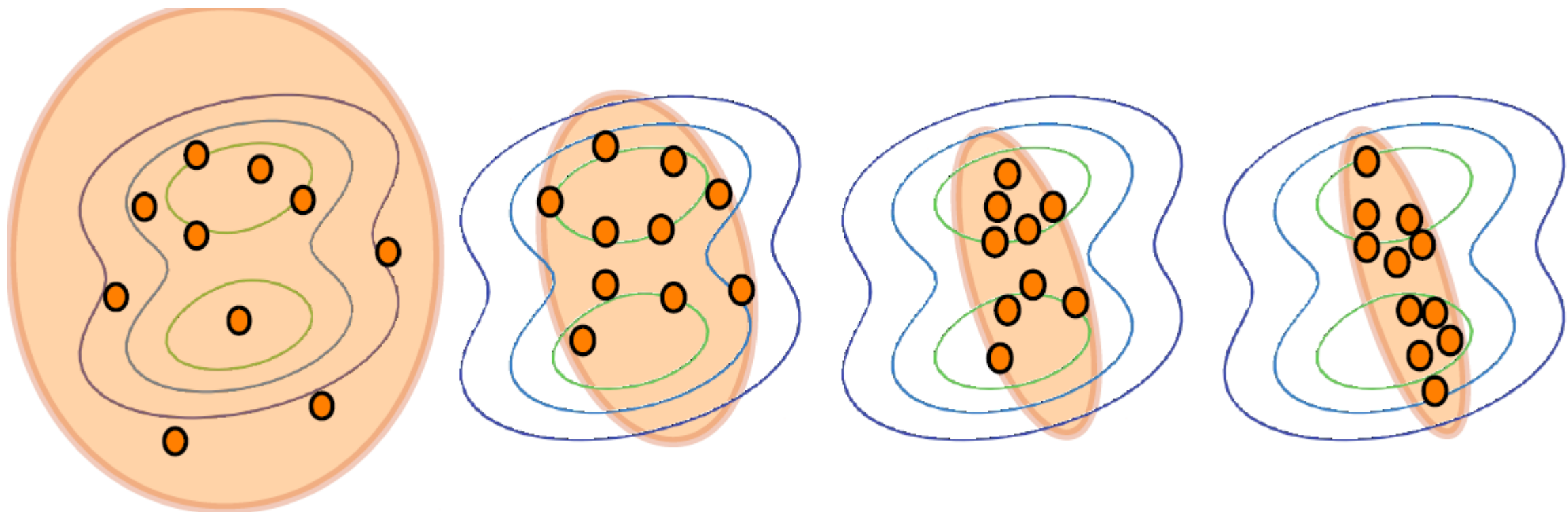
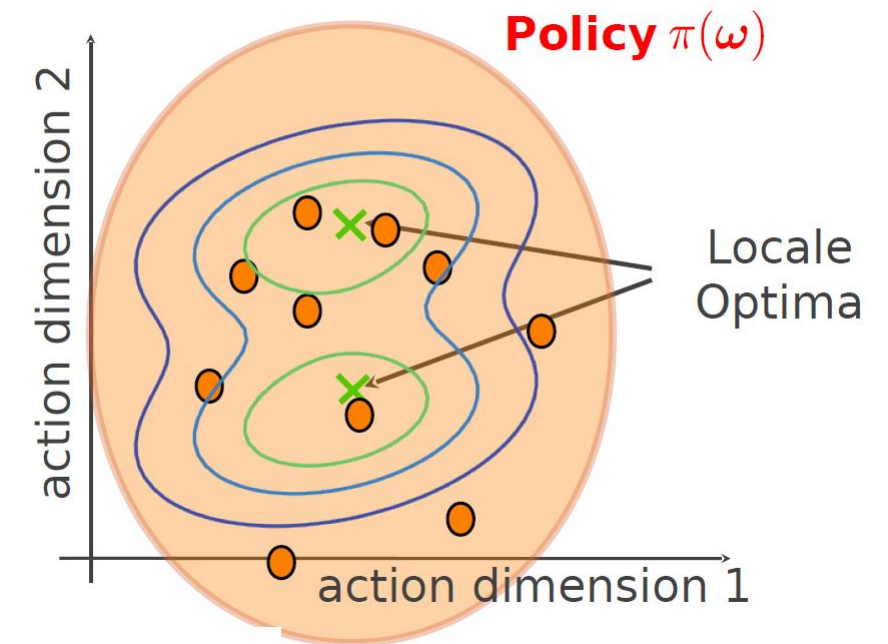
The objective gets „risk attracted“

For moderately stochastic environments it still works well

Illustration on weighted ML



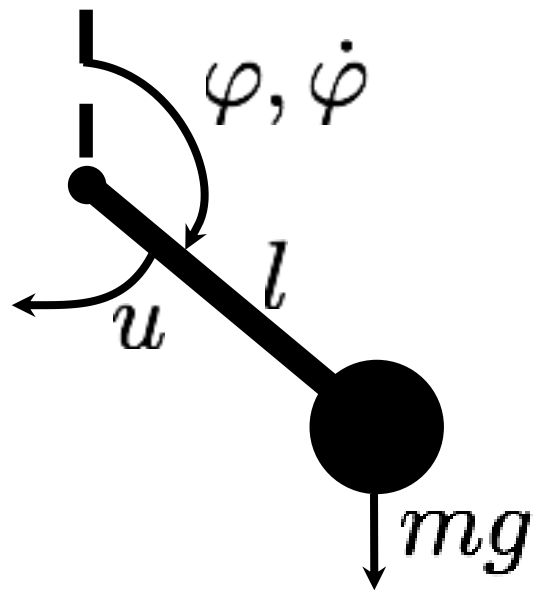
Example for a 2D parameter space:





Underactuated Swing-Up

- swing heavy pendulum up



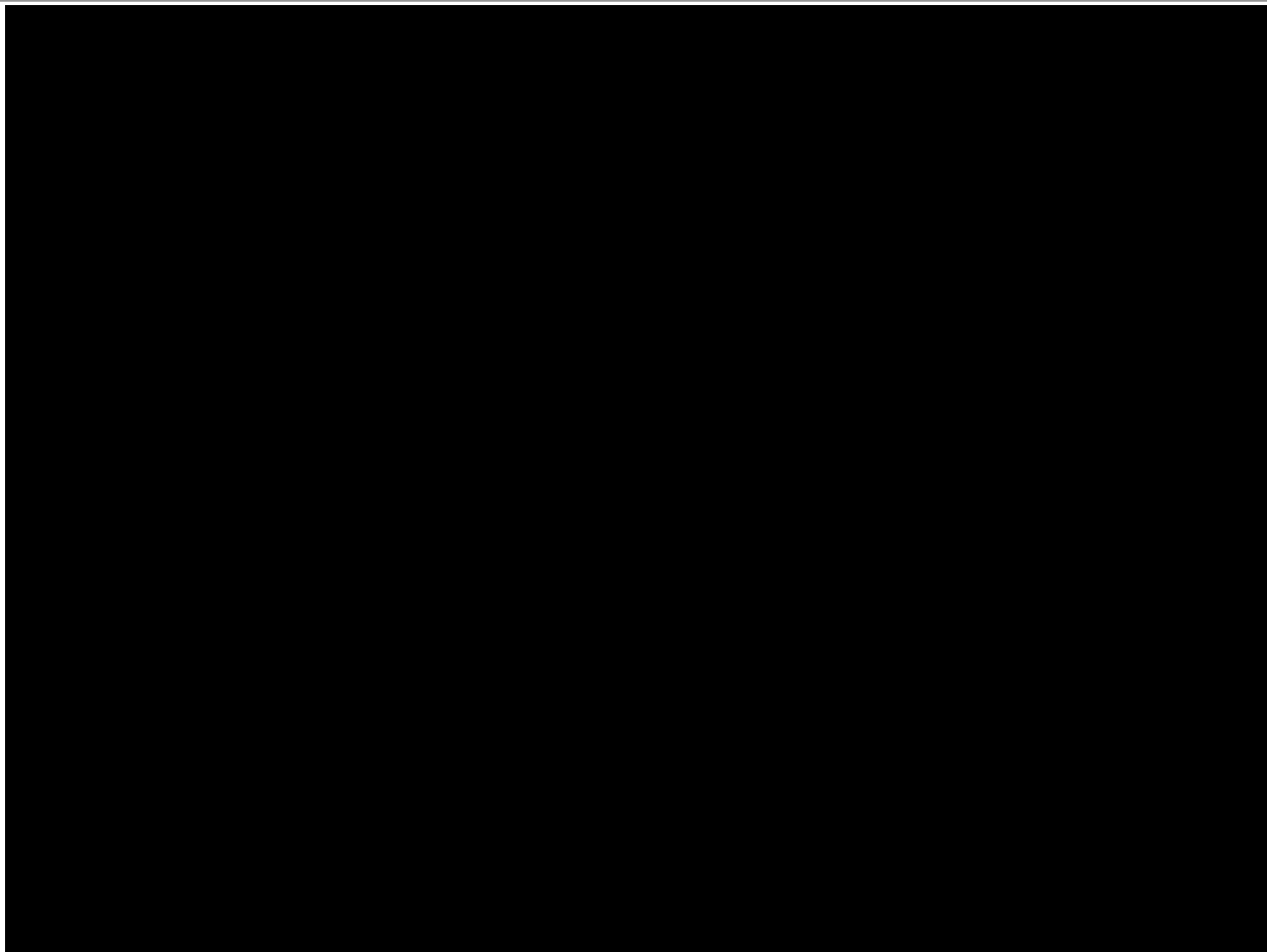
$$ml^2\ddot{\varphi} = -\mu\dot{\varphi} + mgl \sin \varphi + u$$
$$\varphi \in [-\pi, \pi]$$

- motor torques limited, Policy: DMPs

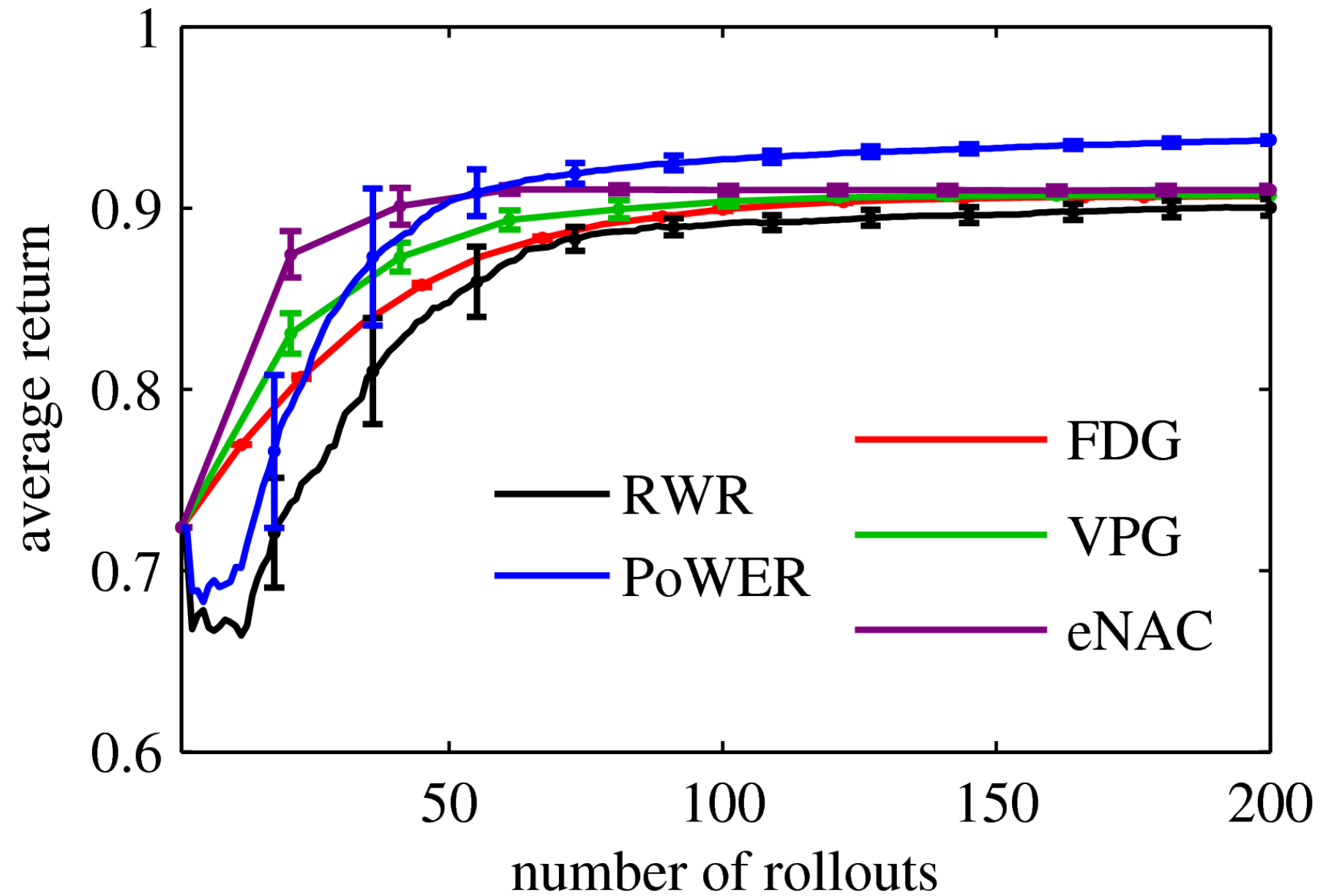
$$|u| \leq u_{max}$$

- reward function

$$r = \exp \left(-\alpha \left(\frac{\varphi}{\pi} \right)^2 - \beta \left(\frac{2}{\pi} \right)^2 \log \cos \left(\frac{\pi}{2} \frac{u}{u_{max}} \right) \right)$$



Underactuated Swing-Up



Ball in the Cup



Video with Ball In the Cup

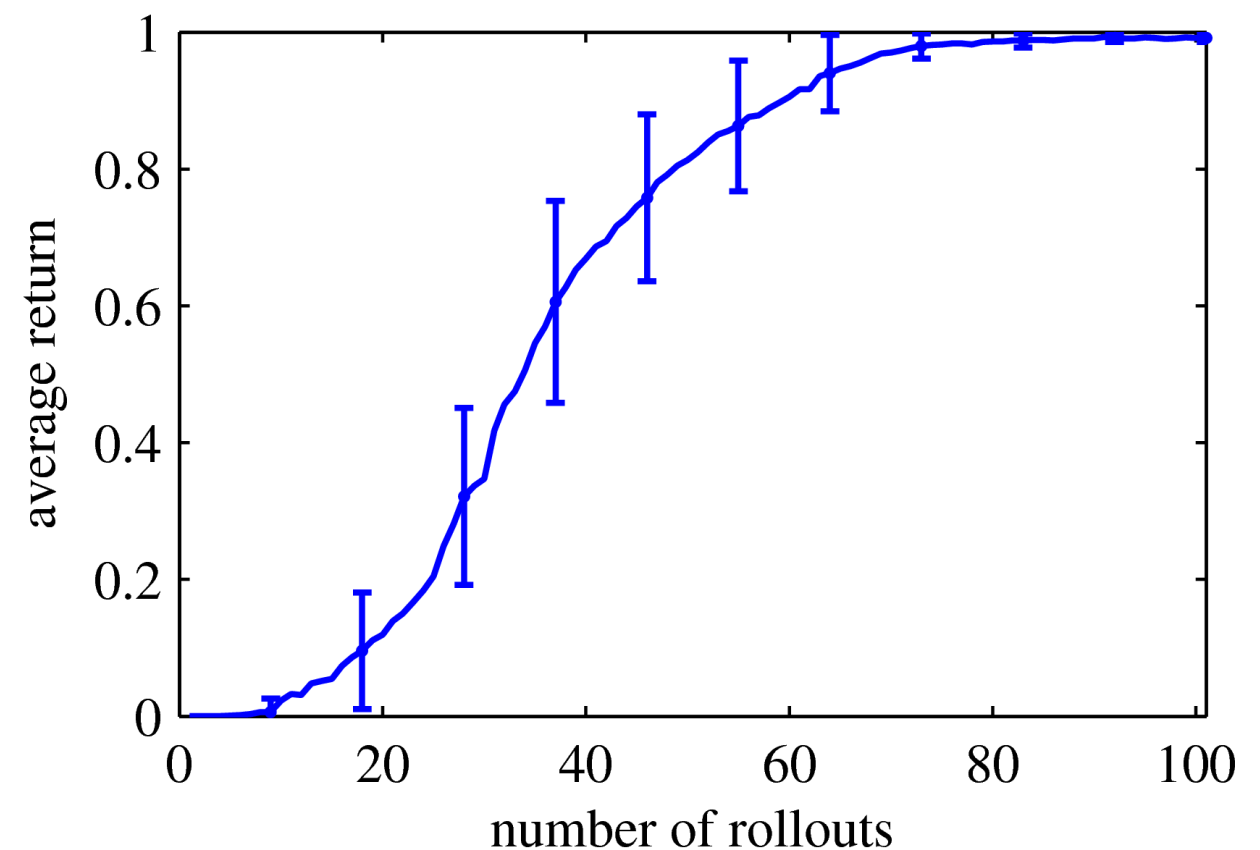
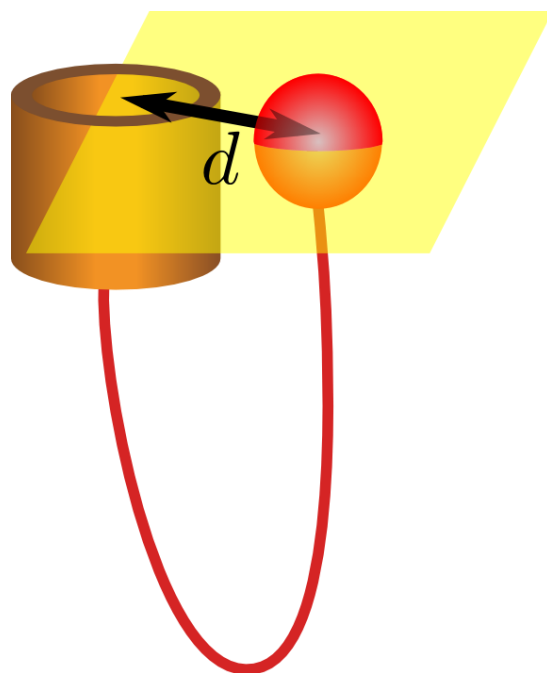
Ball-in-a-Cup



Reward function:

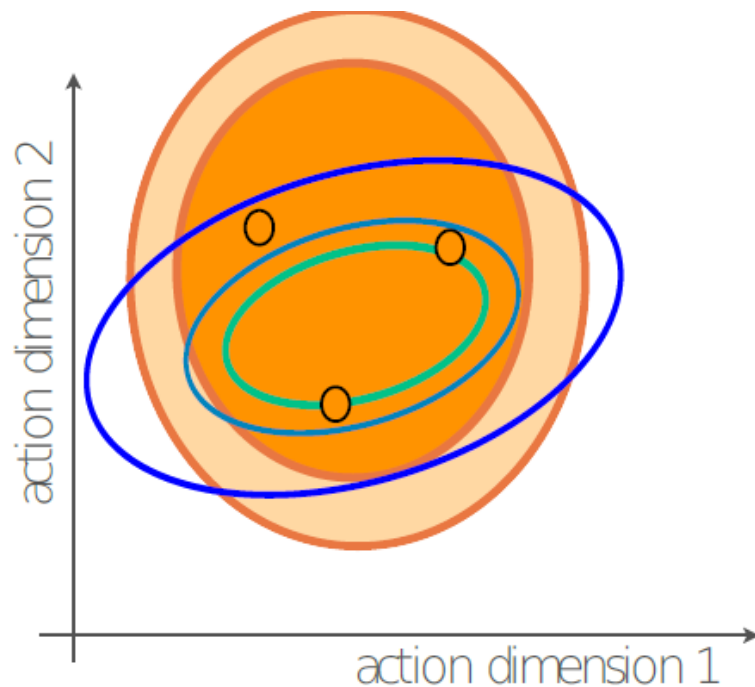
$$r_t = \begin{cases} \exp \left(-\alpha \left((x_c - x_b)^2 + (y_c - y_b)^2 \right) \right) & \text{if } t = t_c \\ 0 & \text{if } t \neq t_c \end{cases}$$

Policy: DMPs

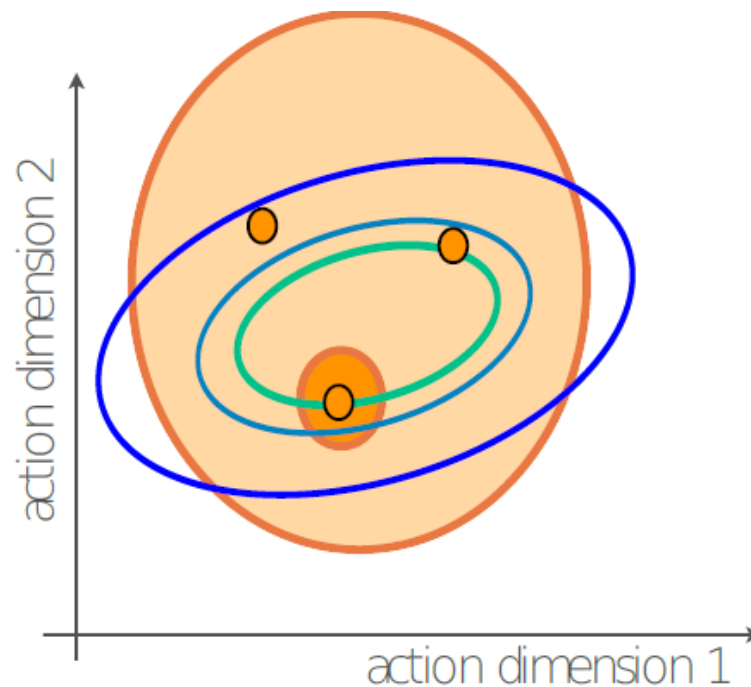


Policy Search: Choosing the metric/step width

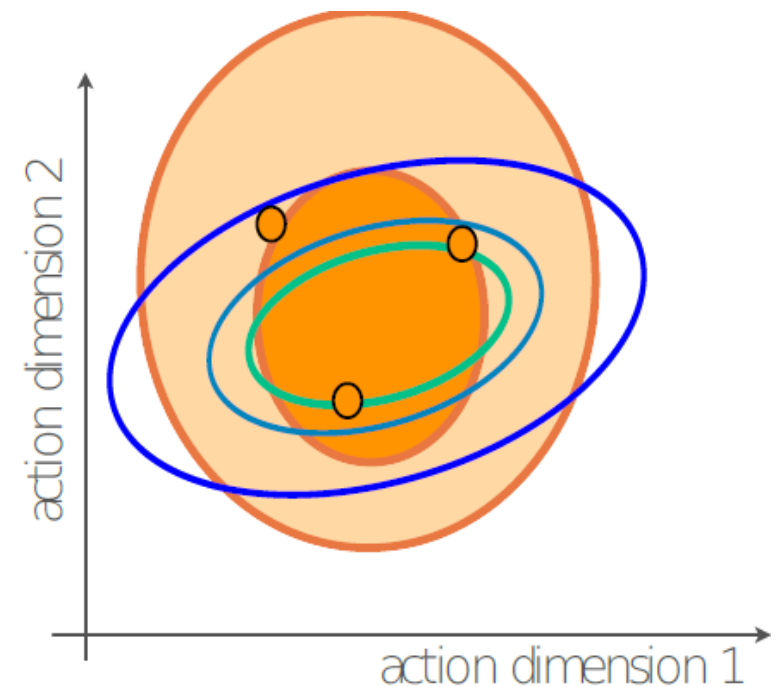
What is a good desired distribution for the policy update?



Too conservative



Too greedy



Moderate tradeoff between exploration and exploitation

We want to have invariance to:

Transformations of the reward

Transformations of the parameter space

For the **weighted ML** methods: How to choose the reward transformation? How to choose the temperature β ?



Outline of the Lecture

1. Introduction
2. Policy Updates by Weighted Maximum Likelihood
3. **Relative Entropy Policy Search (REPS)**
4. REPS for Contextual Policy Search
5. Conclusion

Preserve locality



Coming back to the question: What is a good metric for the policy update?

Goal: Maximize the expected long-term reward

$$J_{\theta} = \mathbb{E}_{\mu_0, \mathcal{P}, \pi} \left[\sum_{t=1}^{T-1} r_t(\mathbf{s}_t, \mathbf{a}_t) + r_T(\mathbf{s}_T) \right]$$

But: We want to preserve locality to achieve a „safe“ policy update

$$\text{s.t.}:: M(\pi_{k+1}, \pi_k) \leq \epsilon$$

M ... metric on the policy update



Used metrics in practice

Euclidian distance in parameter space

$$M(\omega, \omega_{\text{old}}) = \|\omega - \omega_{\text{old}}\|^2$$

➡ Used implicitly by all standard policy gradient approaches

Information-theoretic distance in probability distribution space

$$M(\omega, \omega_{\text{old}}) = \text{KL}(\pi_{\omega} \parallel \pi_{\omega_{\text{old}}})$$

Measures the ‘distance’ between old and new policy in probability space

Policy update: $\text{KL}(\pi_{\omega} \parallel \pi_{\omega_{\text{old}}}) \leq \epsilon$

➡ Invariant to transformations of parameter vector or reward

➡ **Stay close to the data!!**

Used metrics in practice

Two different methods have been used:

Natural Policy Gradient (Peters & Schaal, 2008):

$$\text{KL}(\pi_{\omega} || \pi_{\omega_{\text{old}}}) \approx \Delta\omega^T \mathbf{F} \Delta\omega$$

\mathbf{F} ... Fisher information matrix

Second order Taylor approximation of the KL

 Update is still **done in parameter space**

Relative Entropy Policy Search

Directly **optimize probabilities of the samples** such that

$$\text{KL}(\pi_{\omega} || \pi_{\omega_{\text{old}}}) \leq \epsilon$$

Subsequently, **fit policy to weighted samples** to get the new policy

Relative Entropy Policy Search Peters 2010



$$\max_{\pi} \sum_i \pi(\boldsymbol{\theta}^{[i]}) R(\boldsymbol{\theta}^{[i]})$$

Maximize Reward

s.t:

$$\sum_i \pi(\boldsymbol{\theta}^{[i]}) = 1$$

It's a distribution

$$\text{KL}(\pi(\boldsymbol{\theta}) || q(\boldsymbol{\theta})) \leq \epsilon$$

Stay close to the data

Old Policy

Policy Update is formulated as
constraint optimization problem



Cookbook for Constraint Optimization Problems

Given: **Constraint Optimization Problem**

$$\max_{\mathbf{x}} f(\mathbf{x}) \quad \text{s.t: } g_i(\mathbf{x}) = 0, \quad \forall i$$

Suppose \mathbf{x} is on the valid constraint surface, i.e., $g(\mathbf{x}) = 0$

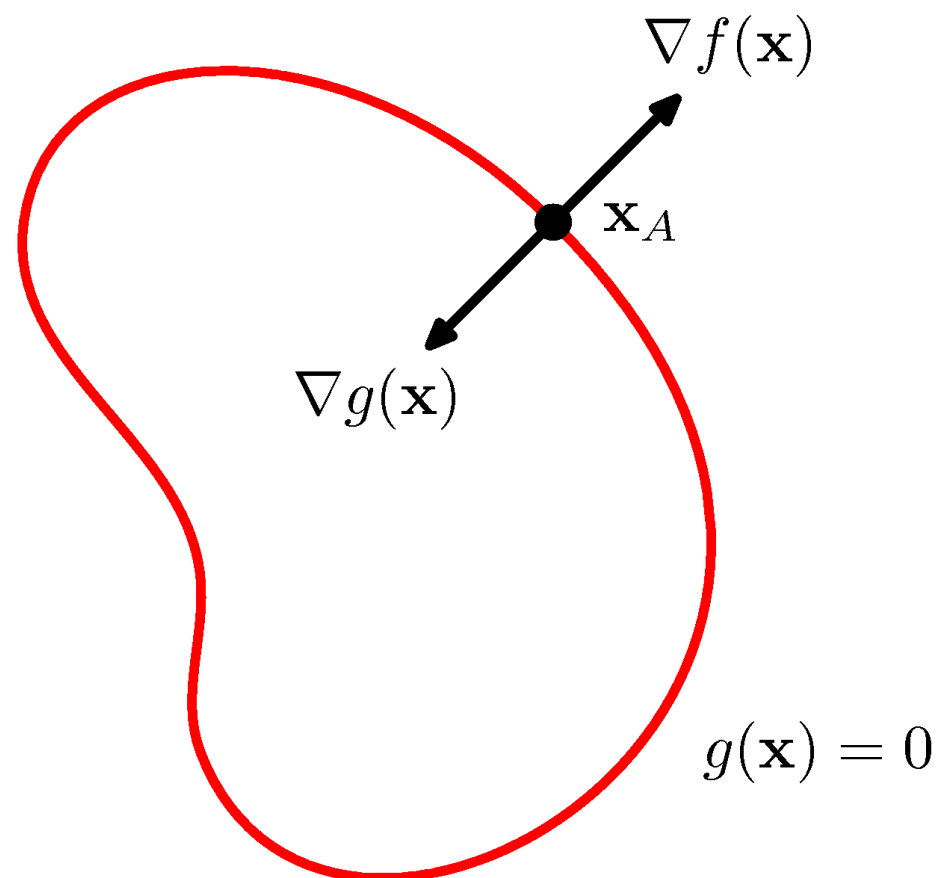
For a Taylor expansion around \mathbf{x}

$$g(\mathbf{x} + \epsilon) \approx g(\mathbf{x}) + \epsilon^T \nabla g(\mathbf{x})$$

From $g(\mathbf{x} + \epsilon) = 0$, it follows that

$$\epsilon^T \nabla g(\mathbf{x}) = 0$$

i.e., $\nabla g(\mathbf{x})$ is normal to the constraint surface





Cookbook for Constraint Optimization Problems

Given: **Constraint Optimization Problem**

$$\max_{\mathbf{x}} f(\mathbf{x}) \quad \text{s.t: } g_i(\mathbf{x}) = 0, \quad \forall i$$

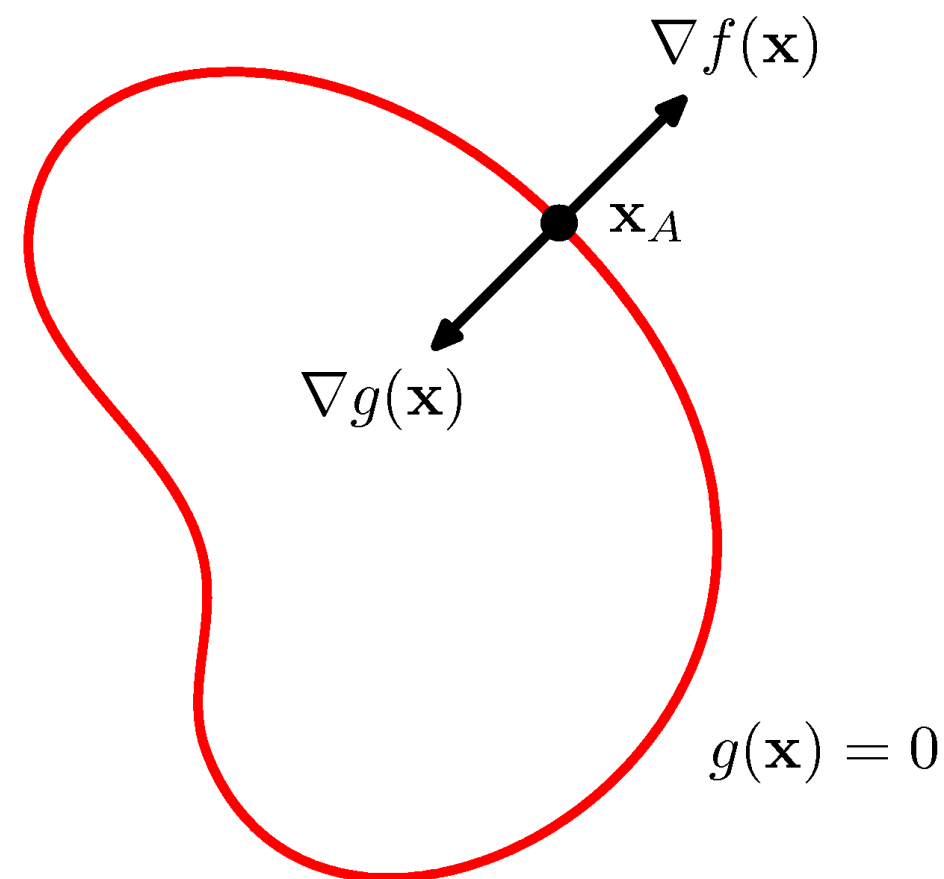
Now we look for a point \mathbf{x}^* on the constraint surface, that maximizes f

$\nabla f(\mathbf{x}^*)$ needs to be orthogonal to the constraint surface (otherwise we could increase f)

⇒ $\nabla f(\mathbf{x}^*)$ and $\nabla g(\mathbf{x})$ are (anti-)parallel

There must be a parameter λ such that

$$\nabla f + \lambda \nabla g = 0$$





Cookbook for Constraint Optimization Problems

Now we can introduce the **Lagrangian function**

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

and we want to find a stationary point of $L(\mathbf{x}, \lambda)$

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda) = 0 \quad \frac{\partial \mathcal{L}(\mathbf{x}, \lambda)}{\partial \lambda} = 0$$

The lagrangian dual function is $g(\boldsymbol{\lambda}) = \max_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$

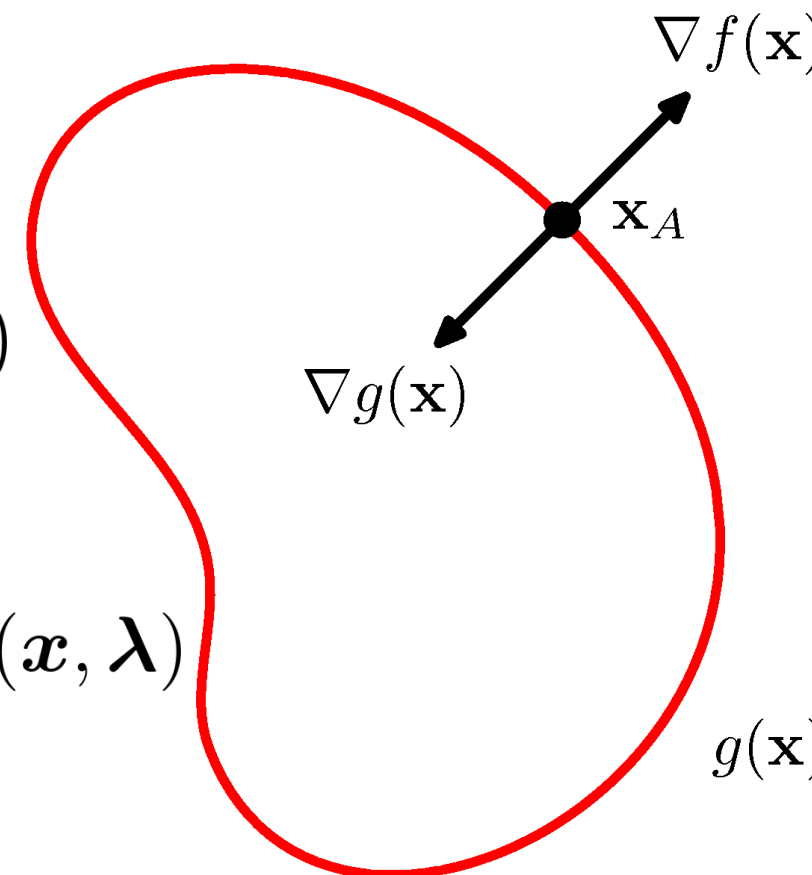
and, hence, $g(\boldsymbol{\lambda}) \geq f(\mathbf{x}) + \lambda g(\mathbf{x})$

For an optimal point \mathbf{x}^* that solves the primal problem

$$g(\boldsymbol{\lambda}) \geq f(\mathbf{x}^*) + \lambda g(\mathbf{x}^*) = f(\mathbf{x}^*)$$

Since $g(\boldsymbol{\lambda}) \geq f(\mathbf{x}^*)$, the optimal λ is obtained by minimizing the dual

$$\lambda^* = \min g(\lambda)$$





Cookbook for Constraint Optimization Problems

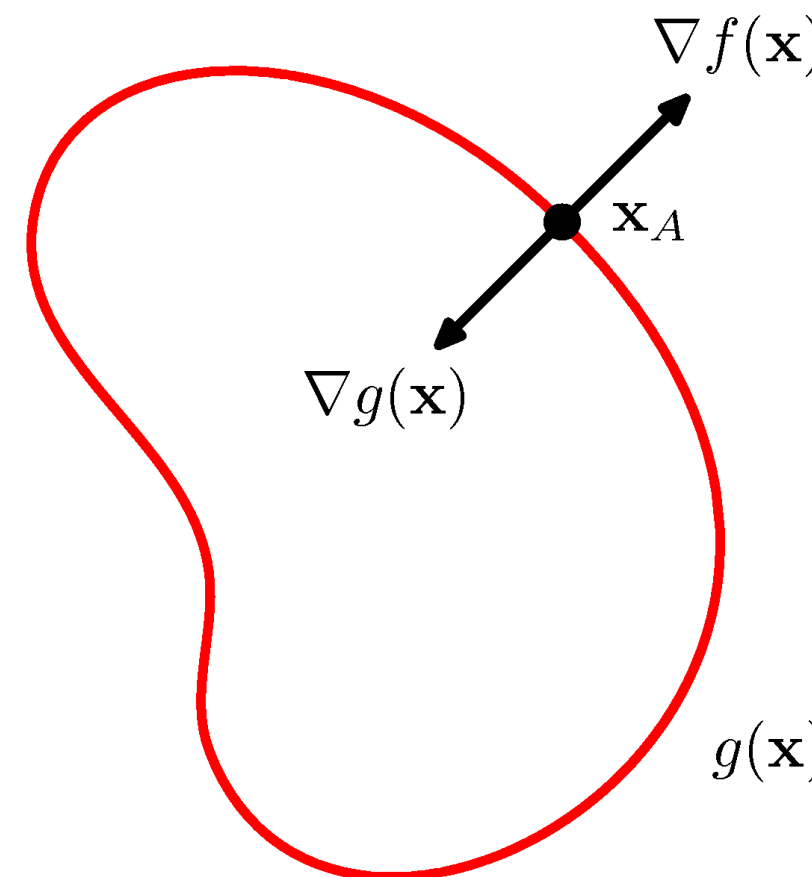
Why should we solve the dual problem?

$$\lambda^* = \min_{\lambda} g(\lambda)$$

It is often easier to solve (**less variables, unconstrained**)

It is **convex**, even if the original problem is not

However, the solution of the dual can be used to solve the **primal only under certain conditions**





Cookbook for Constraint Optimization Problems

Given: Constraint Optimization Problem

$$\max_{\mathbf{x}} f(\mathbf{x}) \quad \text{s.t: } g_i(\mathbf{x}) = 0, \quad \forall i$$

1. Write down **Lagrangian with Lagrangian Multipliers**

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_i \lambda_i g_i(\mathbf{x})$$

2. **Solve for optimal \mathbf{x} for a given $\boldsymbol{\lambda}$**

$$\mathbf{x}^* = \operatorname{argmax}_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$$

3. Set back in Lagrangian to obtain **dual function**

$$g(\boldsymbol{\lambda}) = \max_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$$

4. **Solve for optimal $\boldsymbol{\lambda}^*$: Optimize the (unconstrained) dual function**

$$\boldsymbol{\lambda}^* = \operatorname{argmin}_{\boldsymbol{\lambda}} g(\boldsymbol{\lambda})$$

5. Use $\boldsymbol{\lambda}^*$ to obtain \mathbf{x}^*

This is only a sketch, the theory behind it is more complicated!



Cookbook for Constraint Optimization Problems

Given: Constraint Optimization Problem with several constraints $g_i(x)$

$$\max_x f(x) \quad \text{s.t: } g_i(x) = 0, \quad \forall i$$

1. Write down **Lagrangian with Lagrangian Multipliers**

$$\mathcal{L}(x, \lambda) = f(x) + \sum_i \lambda_i g_i(x)$$

2. **Solve for optimal x for a given λ**

$$x^* = \operatorname{argmax}_x \mathcal{L}(x, \lambda)$$

To see the structure of the solution it is enough to use it until here...

-
3. Set back in Lagrangian to obtain **dual function**

$$g(\lambda) = \sup_x \mathcal{L}(x, \lambda)$$

4. **Solve for optimal λ^* : Optimize the (unconstrained) dual function**

$$\lambda^* = \operatorname{argmin}_\lambda g(\lambda)$$

5. Use λ^* to obtain x^*



REPS: Policy Search as **constraint optimization problem**

$$\max_{\pi} \sum_i \pi(\boldsymbol{\theta}^{[i]}) R(\boldsymbol{\theta}^{[i]})$$

Maximize Reward

s.t: $\sum_i \pi(\boldsymbol{\theta}^{[i]}) = 1$

It's a distribution

$$\text{KL}(\pi(\boldsymbol{\theta}) || q(\boldsymbol{\theta})) \leq \epsilon$$

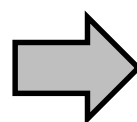
Stay close to the data

Scalingfactor η

- **Automatically chosen from optimization** (Lagrange Multiplier)
- Specified by KL-bound ϵ
- We get the exponential transformation (used by EM) for free

$$\pi(\boldsymbol{\theta}) \propto q(\boldsymbol{\theta}) \exp \left(\frac{\mathcal{R}_{\boldsymbol{\theta}}}{\eta} \right)$$

Analytic Solution



Weight computed by REPS

$$w^{[i]} \propto \exp \left(\frac{\mathcal{R}^{[i]}}{\eta} \right)$$



Getting the Lagrangian multipliers

How to get η :

Solve **dual optimization** problem:

Dual function:
$$g(\eta) = \eta\epsilon + \eta \log \int q(\boldsymbol{\theta}) \exp \left(\frac{\mathcal{R}\boldsymbol{\theta}}{\eta} \right) d\boldsymbol{\omega}$$

$$= \eta\epsilon + \eta \log \sum_{i=1}^N \frac{1}{N} \exp \left(\frac{\mathcal{R}^{[i]}}{\eta} \right)$$

Minimize: $\eta^* = \operatorname{argmin}_{\eta} g(\eta) \quad \text{s.t: } \eta > 0$

Log-sum-exp softmax structure

Optimized by standard optimization tools

(e.g. trust region algorithms)

Episode-based policy evaluation strategy



Contextual Policy Search

Context \mathbf{x} describes objectives of the task (fixed before task execution)

E.g.: Target location to throw a ball

We now want to learn an upper level policy that adapts θ to the context $\pi(\theta|\mathbf{x};\omega)$

Data-set used for policy update

$$\mathcal{D}_{\text{episode}} = \{\theta^{[i]}, \mathbf{x}^{[i]}, R^{[i]}\}_{i=1\dots N}$$

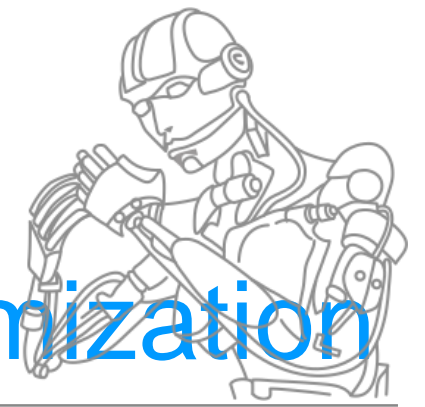
Goal: maximize expected reward

$$J_{\pi} = \iint \mu_0(\mathbf{x}) \pi(\theta|\mathbf{x}) \mathcal{R}_{\mathbf{x}\theta} d\mathbf{x} d\theta$$



Outline of the Lecture

1. Introduction
2. Policy Updates by Weighted Maximum Likelihood
3. Relative Entropy Policy Search (REPS)
4. **REPS for Contextual Policy Search**
5. Conclusion



Contextual Policy Search as **constraint optimization**

We now **optimize over the joint distribution** $p(\mathbf{x}, \boldsymbol{\theta}) = \mu(\mathbf{x})\pi(\boldsymbol{\theta}|\mathbf{x})$ in order to be able to use context-parameter pairs instead of many parameters for a single context

$$\max_p \sum_i p(\mathbf{x}^{[i]}, \boldsymbol{\theta}^{[i]}) R(\mathbf{x}^{[i]}, \boldsymbol{\theta}^{[i]})$$

Maximize Reward

$$\sum_i p(\mathbf{x}^{[i]}, \boldsymbol{\theta}^{[i]}) = 1$$

It's a distribution

$$\text{KL}(p(\mathbf{x}, \boldsymbol{\theta}) || q(\mathbf{x}, \boldsymbol{\theta})) \leq \epsilon$$

Stay close to the data

$$p(\mathbf{x}) = \sum_{\boldsymbol{\theta}} p(\mathbf{x}, \boldsymbol{\theta}) = \mu_0(\mathbf{x})$$

Reproduce given context distribution $\mu_0(\mathbf{x})$

Continuous Context?



Adding the context constraints

$$\max_p \sum_i p(\mathbf{x}^{[i]}, \boldsymbol{\theta}^{[i]}) R(\mathbf{x}^{[i]}, \boldsymbol{\theta}^{[i]})$$

Maximize Reward

$$\sum_i p(\mathbf{x}^{[i]}, \boldsymbol{\theta}^{[i]}) = 1$$

It's a distribution

$$\text{KL}(\pi(\boldsymbol{\theta}|\mathbf{x})\mu(\mathbf{x})||q(\mathbf{x}, \boldsymbol{\theta})) \leq \epsilon$$

Stay close to the data

$$\sum_{\mathbf{x}} p(\mathbf{x}) \phi(\mathbf{x}) = \hat{\phi}$$

Reproduce given context
feature averages

$$\text{e.g., } \phi(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$$



Match Mean and Variance

$$\mu(\mathbf{x})\pi(\boldsymbol{\theta}|\mathbf{x}) \propto q(\mathbf{x}, \boldsymbol{\theta}) \exp\left(\frac{R_{\mathbf{x}\boldsymbol{\theta}} - V(\mathbf{x})}{\eta}\right)$$



Adding the context constraints

$$\max_{\pi, \mu} \sum_i \mu(\mathbf{x}^{[i]}) \pi(\boldsymbol{\theta}^{[i]} | \mathbf{x}^{[i]}) R(\mathbf{x}^{[i]}, \boldsymbol{\theta}^{[i]})$$

Maximize Reward

$$\sum_i \pi(\boldsymbol{\theta}^{[i]} | \mathbf{x}^{[i]}) \mu(\mathbf{x}^{[i]}) = 1$$

It's a distribution

$$KL(\pi(\boldsymbol{\theta} | \mathbf{x}) \| \pi(\boldsymbol{\theta})) \leq \epsilon$$

Stay close to the prior

Baseline $V(\mathbf{x}) = \phi^T(\mathbf{x})\mathbf{v}$

Context-dependent reward baseline, linear in features

\mathbf{v} is specified by feature constraint, Lagrange Multiplier

e.g., $\phi(\mathbf{x}) = \begin{bmatrix} x^2 \end{bmatrix}$

Match mean and variance

$$p(\mathbf{x}, \boldsymbol{\theta}) \propto q(\mathbf{x}, \boldsymbol{\theta}) \exp \left(\frac{R_{\mathbf{x}\boldsymbol{\theta}} - V(\mathbf{x})}{\eta} \right)$$



Getting the Lagrangian multipliers

How to get η, \mathbf{v}

Solve **dual optimization** problem:

Dual function:

$$g(\eta, \mathbf{v}) = \eta\epsilon + \hat{\phi}^T \mathbf{v} + \eta \log \sum_i \frac{1}{N} \exp \left(\frac{\mathcal{R}^{[i]} - \phi^T(\mathbf{x}^{[i]})\mathbf{v}}{\eta} \right)$$

Minimize: $[\eta^*, \mathbf{v}^*] = \operatorname{argmin}_{\eta} g(\eta, \mathbf{v}) \quad \text{s.t: } \eta > 0$

Integral is over the context-parameter space

We can use $(\mathbf{x}^{[i]}, \boldsymbol{\theta}^{[i]})$ samples instead of many samples $\boldsymbol{\theta}^{[ij]}$ per context $\mathbf{x}^{[i]}$



Policy generalization with weighted ML

Estimate parametric policy $\pi_{\omega}(\theta|x)$:

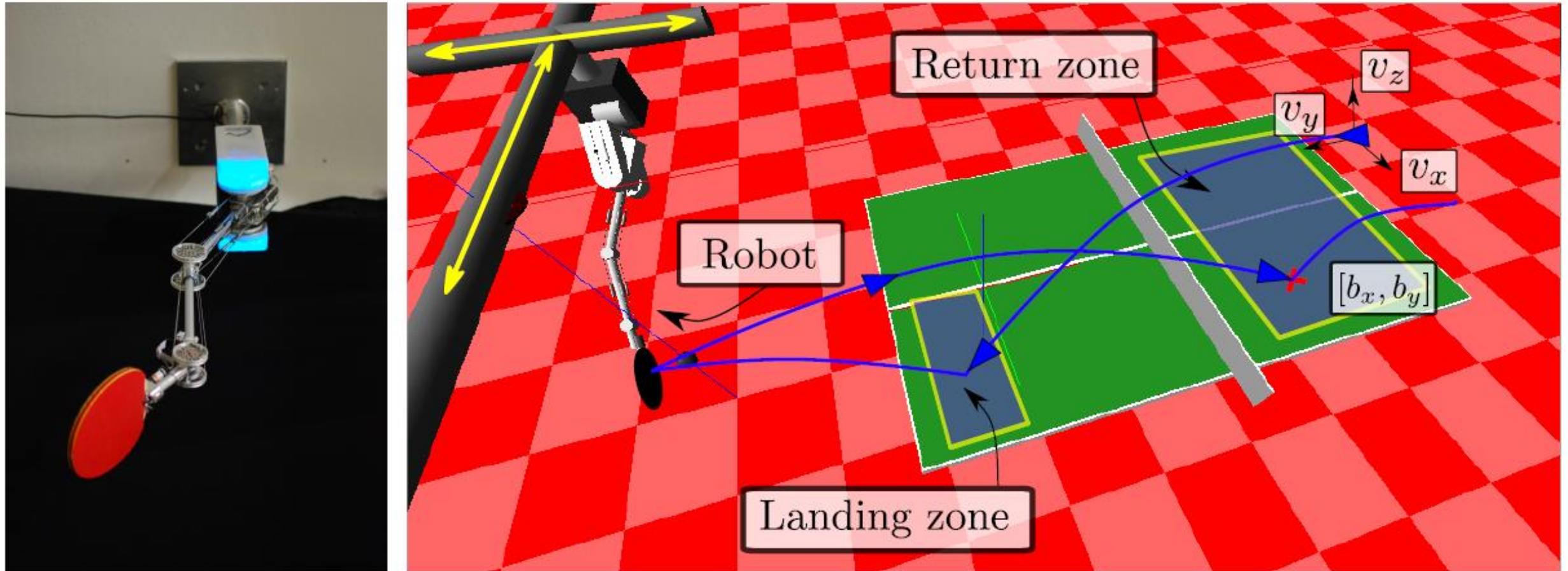
If $\pi_{\omega}(\theta|x) = \mathcal{N}(\theta|Kx + k, \Sigma_{\omega})$ is Gaussian:

$$\begin{bmatrix} k^T \\ K^T \end{bmatrix} = (\underbrace{X^T}_{\text{Input Matrix}} \underbrace{D}_{\text{Weighting Matrix}} X)^{-1} X^T \underbrace{DA}_{\text{Parameter Matrix}} \mu_i)^T$$

$$\mu_i = k + Kx_i \quad \underline{\mu} = \frac{\sum_i p_i \mu_i)^T}{\sum_i p_i}$$

Just standard **weighted linear regression**...

Table tennis experiments



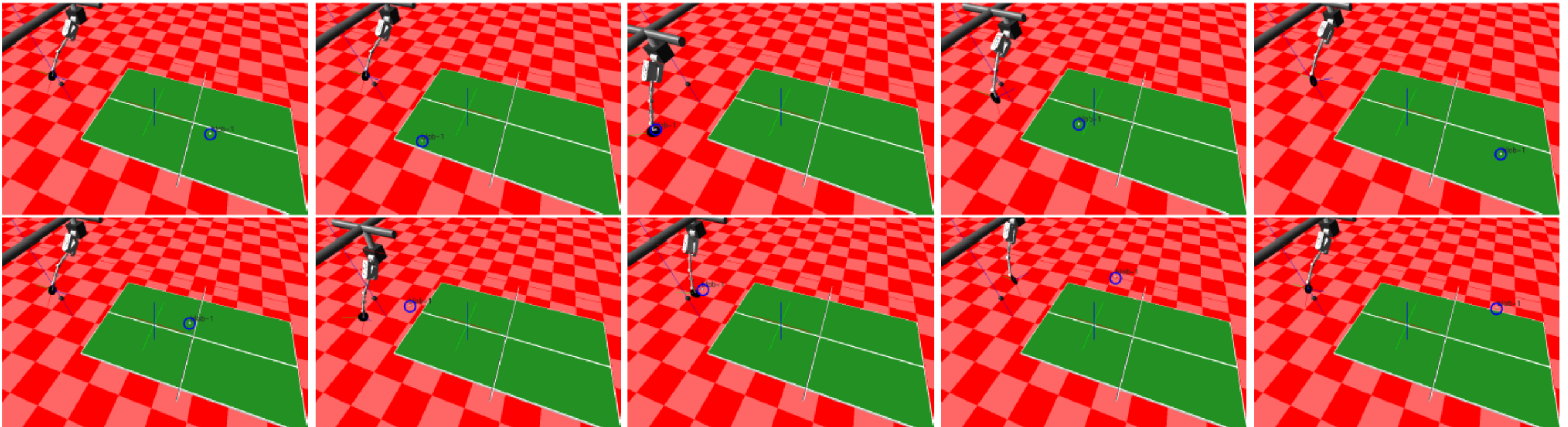
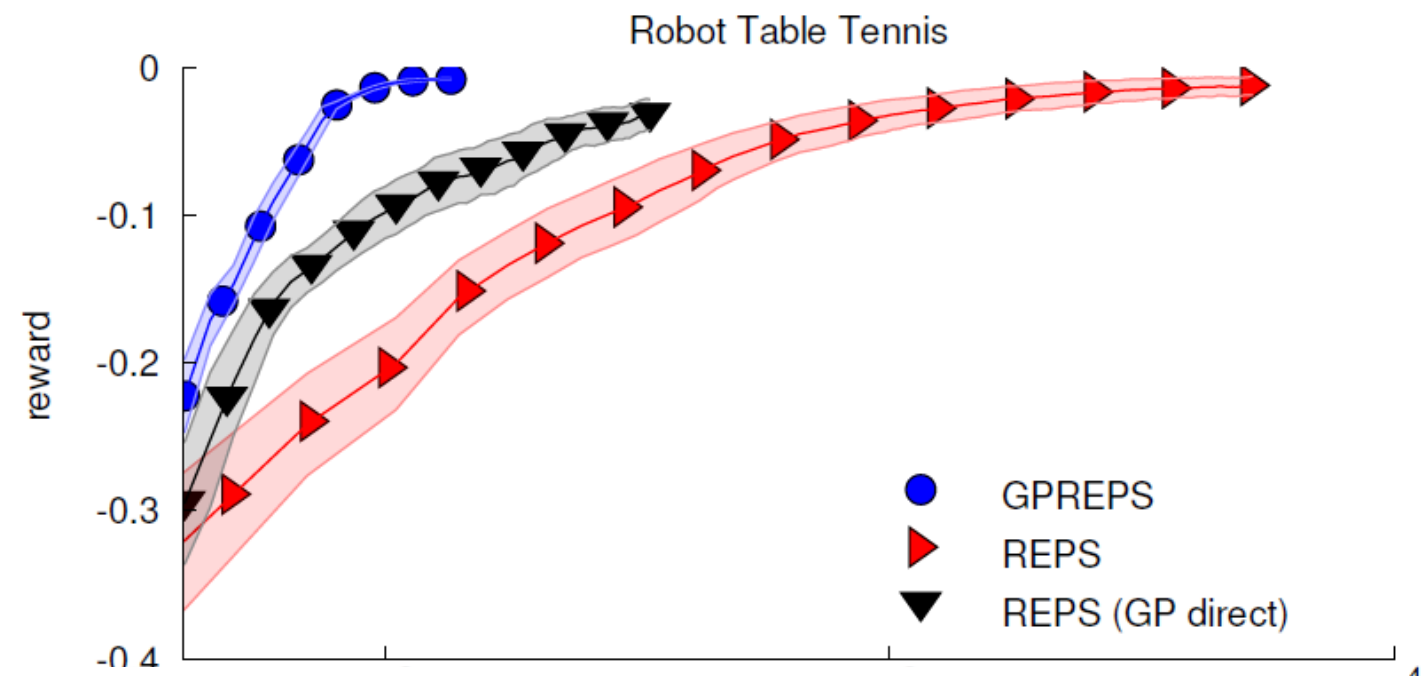
[Kupscik, Neumann et al, submitted, 2013]

Table tennis experiments



REPS with learned forward models

- Complex behavior can be learned within 100 episodes



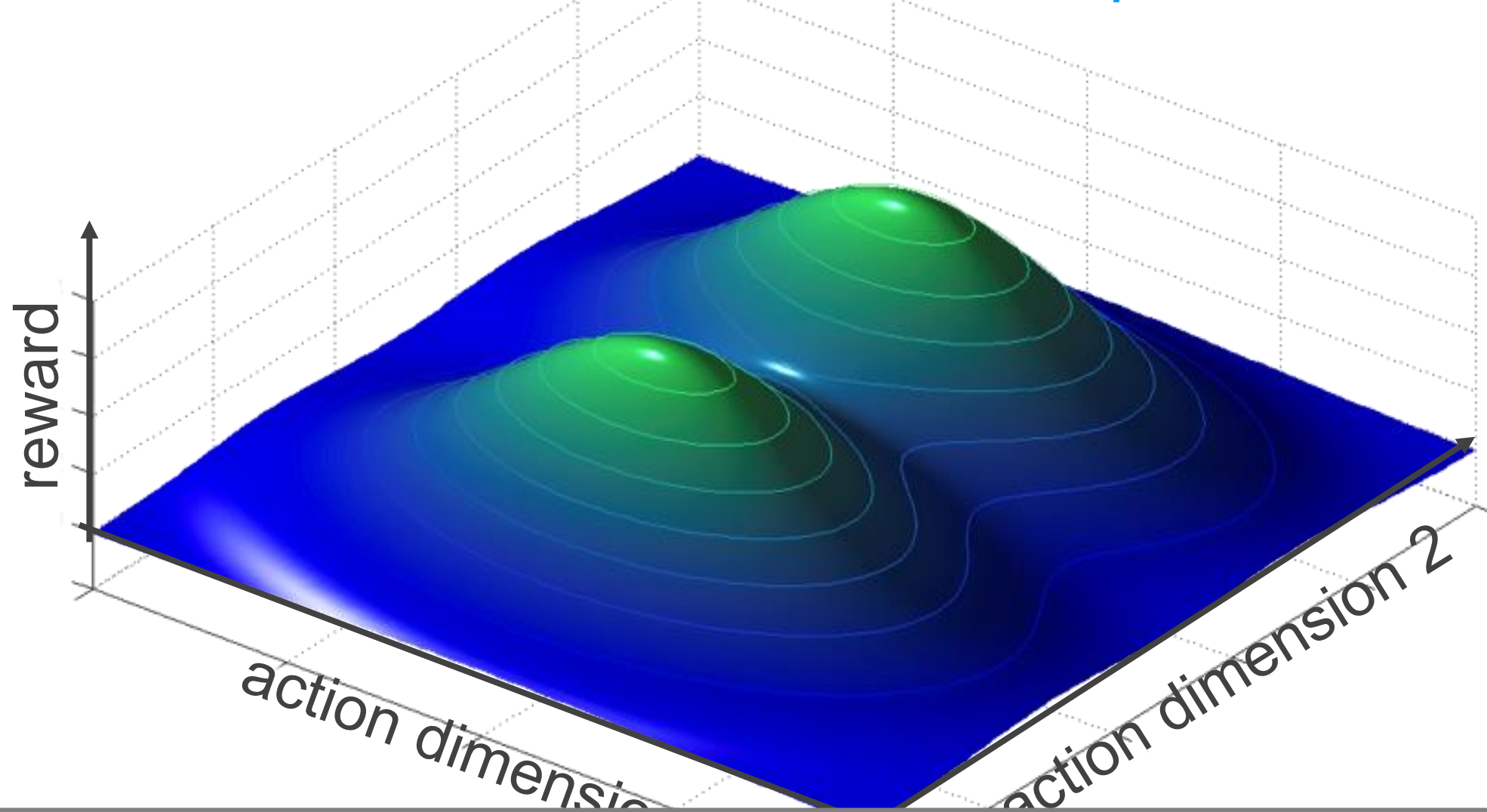


Outline of the Lecture

1. Introduction
2. Policy Updates by Weighted Maximum Likelihood
3. Relative Entropy Policy Search (REPS)
4. REPS for Contextual Policy Search
5. Learning versatile solutions
6. Conclusion

Versatile Solutions: Illustration

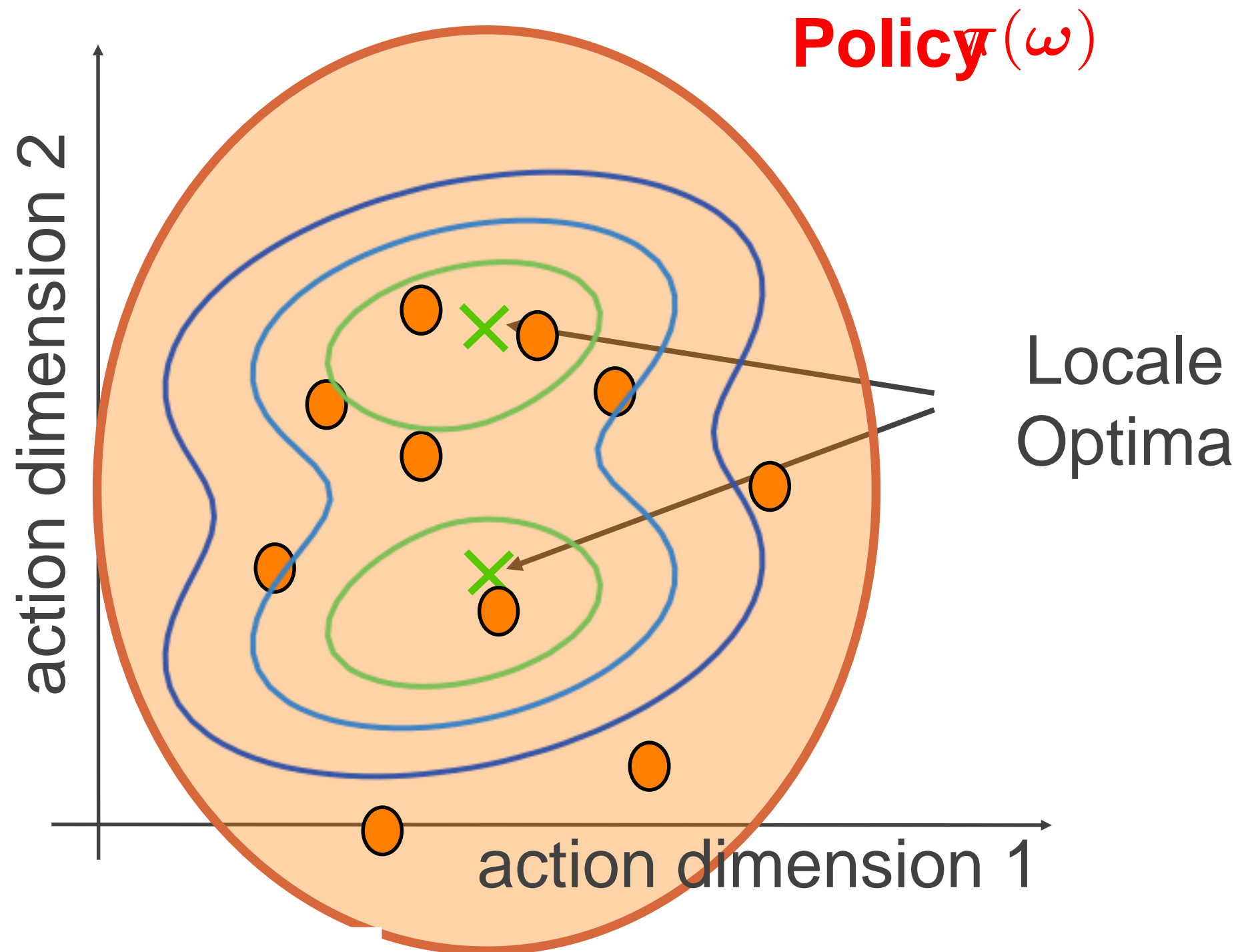
Motor-tasks have **multiple solutions**



C. Daniel, G. Neumann, J. Peters, *Hierarchical Relative Entropy Policy Search*, AISTATS 2012

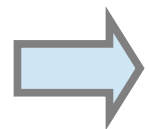
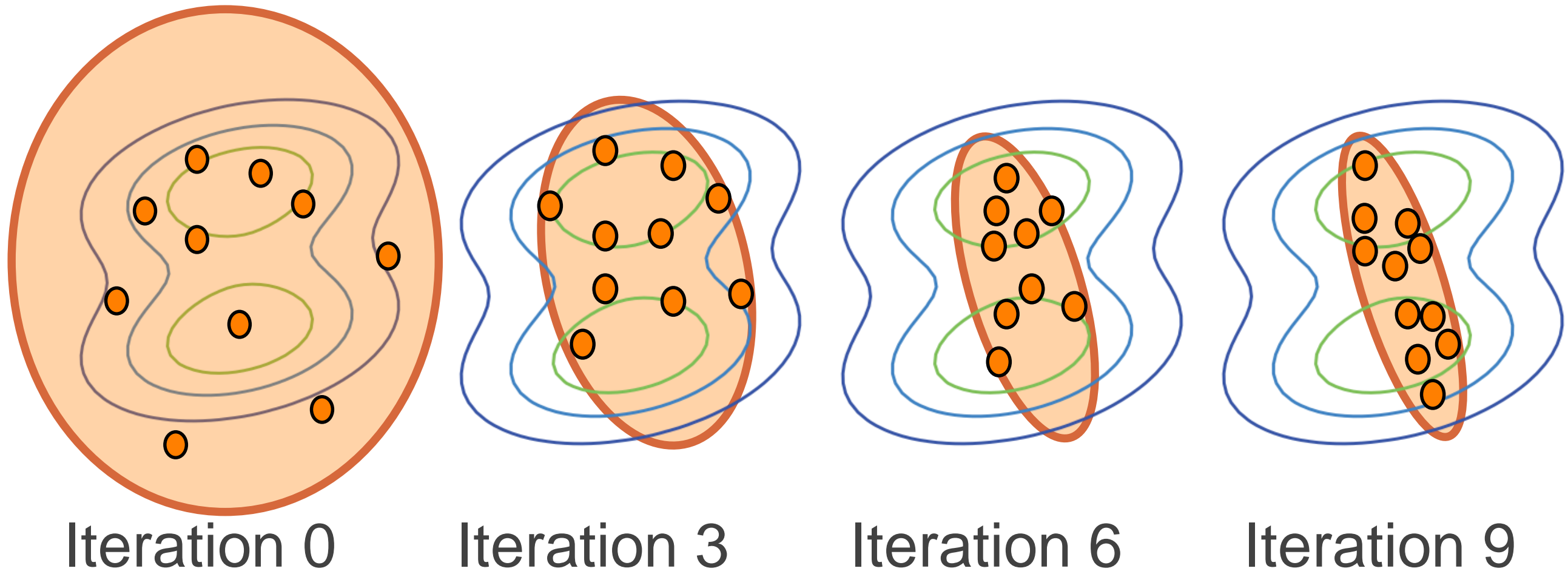
C. Daniel, G. Neumann, J. Peters, *Learning Concurrent Motor Skills in Versatile Solution Spaces*, IROS 2012,
Best Cognitive Systems Paper, Best Paper Finalist

Illustration



Illustration

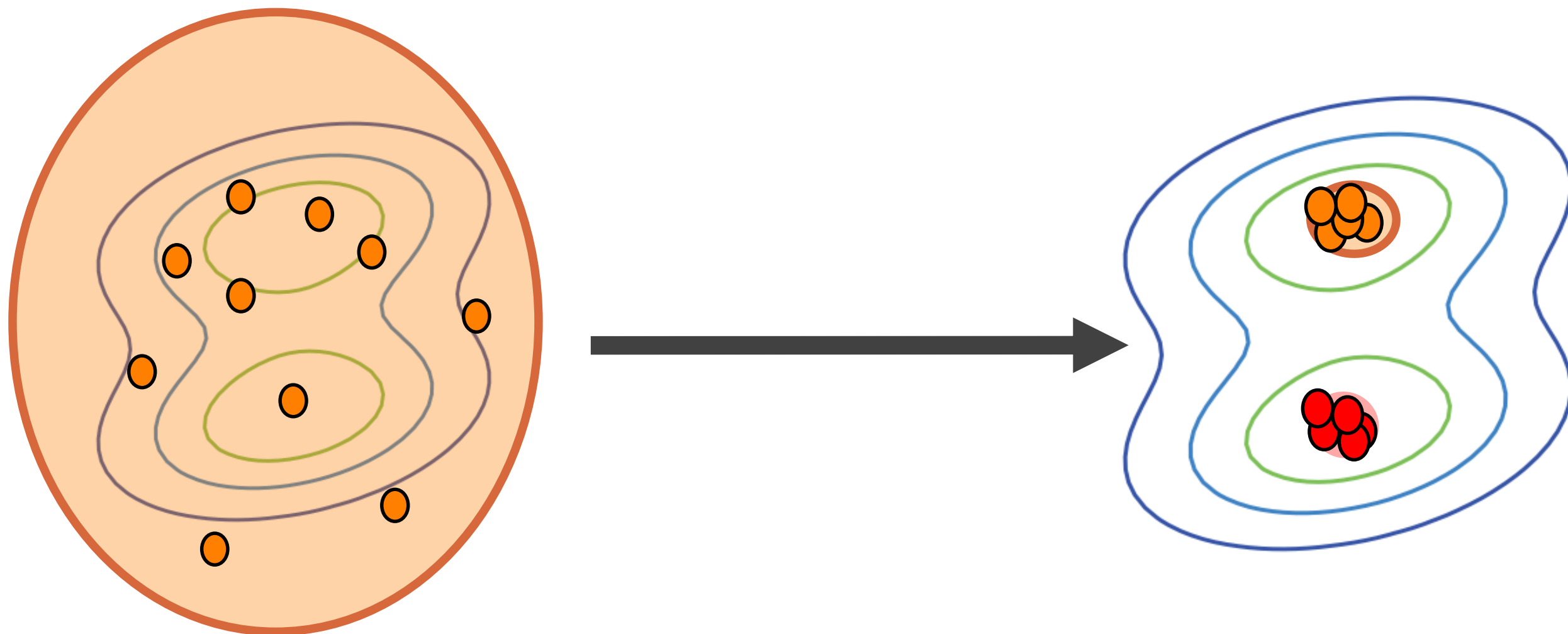
Current Formulation:



Policy averages over several modes.

Illustration

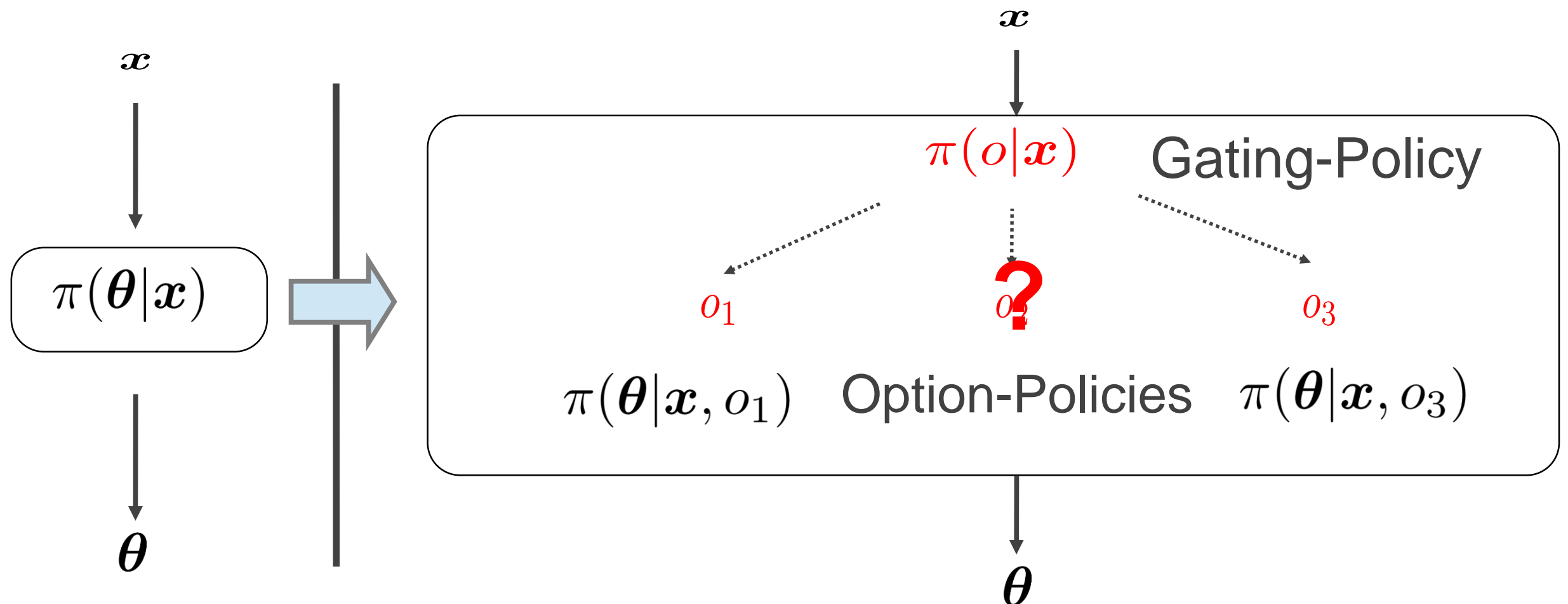
We want to find **both** solutions!



Introduce Hierarchy

Upper-level policy as combination of options

- Selection of the option: **Gating-policy**
- Selection of the parameters: **Option-policy**



“Naive” Hierarchical Approach

$$\max_{\pi, \mu} \sum_{\mathbf{x}, \omega, \mathbf{o}} \mu(\mathbf{x}) \pi(\omega | \mathbf{x}, \mathbf{o}) \pi(\mathbf{o} | \mathbf{x}) R_{\mathbf{x}\omega}$$

Maximize reward

$$\sum_{\mathbf{x}, \omega, \mathbf{o}} \mu(\mathbf{x}) \pi(\omega | \mathbf{x}, \mathbf{o}) \pi(\mathbf{o} | \mathbf{x}) = 1$$

Distribution

$$\sum_{\mathbf{x}} \mu(\mathbf{x}) \phi(\mathbf{x}) = \hat{\phi}$$

Reproduce Context-Features

$$\epsilon \geq \sum_{\mathbf{x}, \omega, \mathbf{o}} \mu(\mathbf{x}) \pi(\omega, \mathbf{o} | \mathbf{x}) \log \frac{\mu(\mathbf{x}) \pi(\omega, \mathbf{o} | \mathbf{x})}{q(\mathbf{x}, \omega, \mathbf{o})}$$

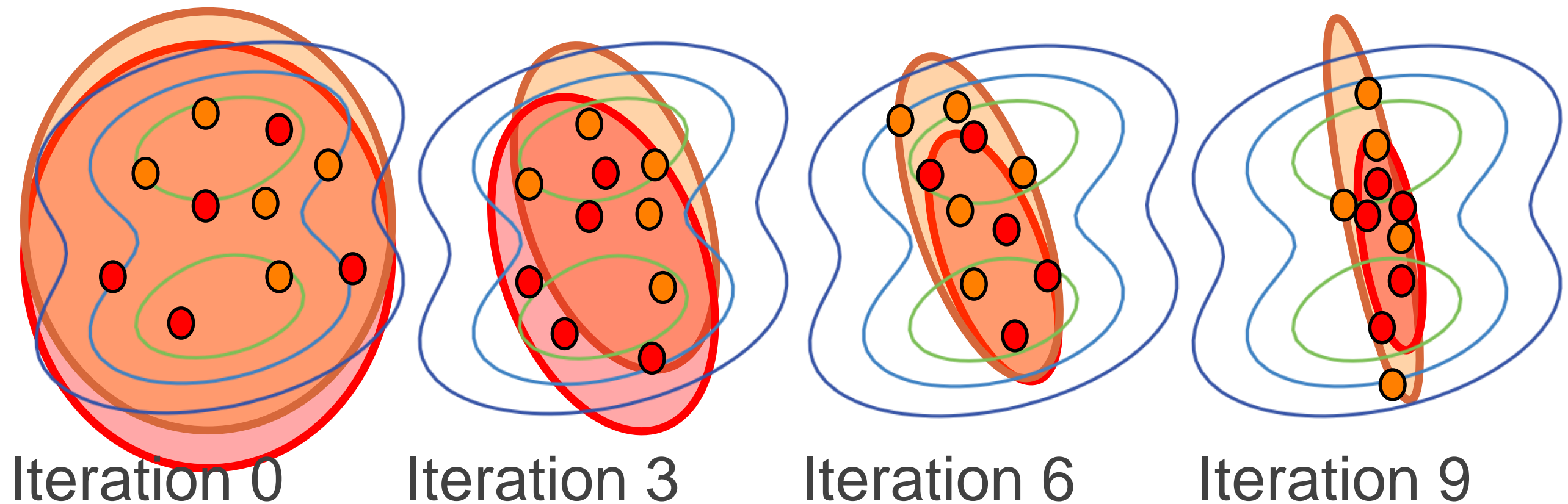
Stay close to the “data”

?

Versatile Solutions

Illustration

“Naive” Approach:

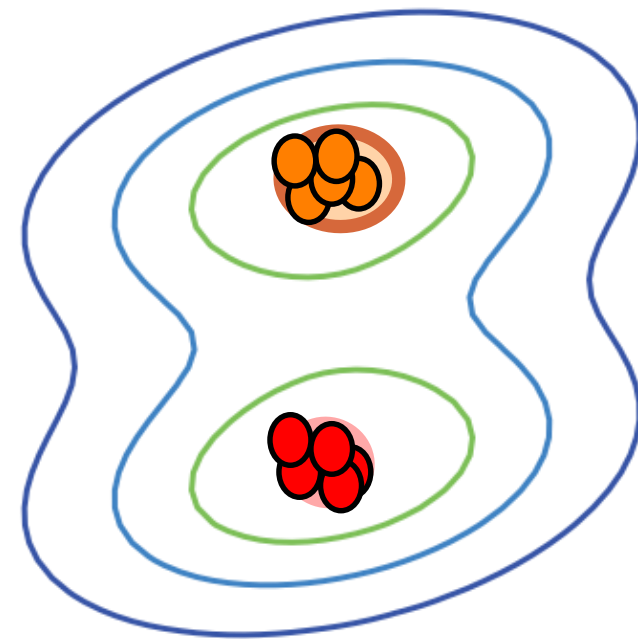


Multiple Options, **BUT** no separation

Learning versatile Options

Options should represent distinct solutions.

➡ Limit the overlap of the options



High entropy of $p(o|\mathbf{x}, \boldsymbol{\theta})$ ➡ overlap

Limit the entropy ➡ less overlap

$$\kappa \geq \mathbb{E} \left[\underbrace{- \sum_o p(o|\mathbf{x}, \boldsymbol{\theta}) \log p(o|\mathbf{x}, \boldsymbol{\theta})}_{\text{Entropy}} \right]$$

Hierarchical REPS (HiREPS)

$$\max_{\pi, \mu} \sum_{\mathbf{x}, \omega, \mathbf{o}} \mu(\mathbf{x}) \pi(\omega | \mathbf{x}, \mathbf{o}) \pi(\mathbf{o} | \mathbf{x}) R_{\mathbf{x}\omega}$$

$$\sum_{\mathbf{x}, \omega, \mathbf{o}} \mu(\mathbf{x}) \pi(\omega | \mathbf{x}, \mathbf{o}) \pi(\mathbf{o} | \mathbf{x}) = 1$$

$$\sum_{\mathbf{x}} \mu(\mathbf{x}) \phi(\mathbf{x}) = \hat{\phi}$$

$$\epsilon \geq \sum_{\mathbf{x}, \omega, \mathbf{o}} \mu(\mathbf{x}) \pi(\omega, \mathbf{o} | \mathbf{x}) \log \frac{\mu(\mathbf{x}) \pi(\omega, \mathbf{o} | \mathbf{x})}{q(\mathbf{x}, \omega, \mathbf{o})}$$

Maximize reward

Distribution

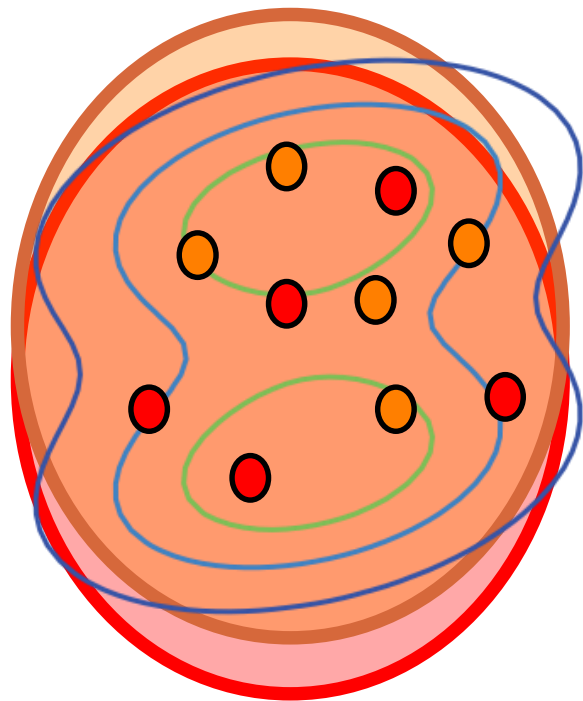
Reproduce Context-Features

Stay close to the “data”, no wild exploration

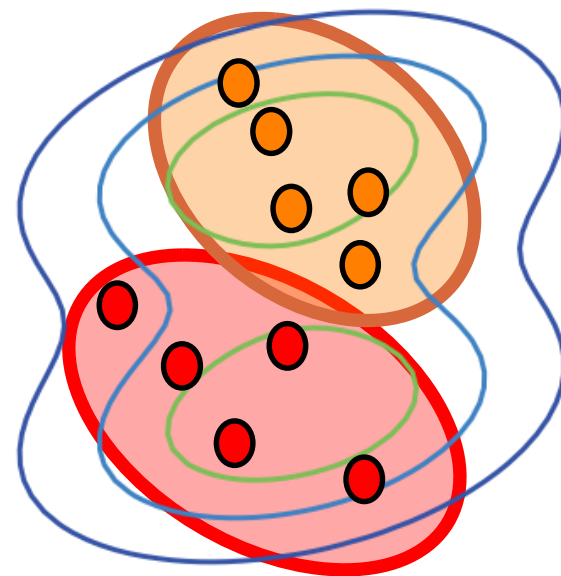
$$\kappa \geq \mathbb{E} [-p(\mathbf{o} | \mathbf{x}, \omega) \log p(\mathbf{o} | \mathbf{x}, \omega)]$$

Versatile Solutions

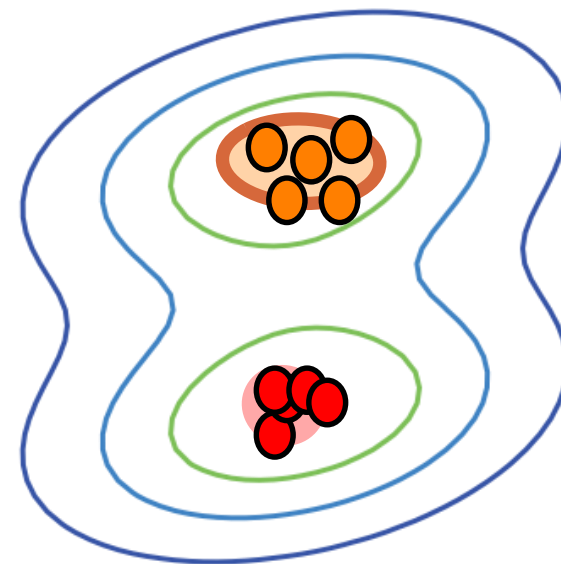
HiREPS



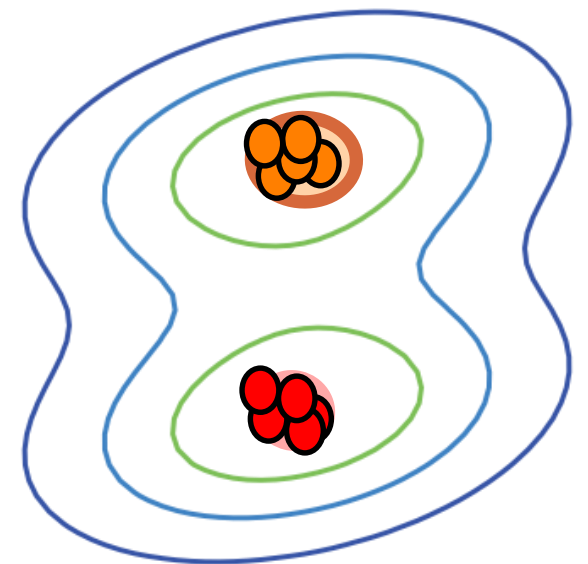
Iteration 0



Iteration 3



Iteration 6



Iteration 9

Learning of versatile, distinct solutions due to separation of options.

Tetherball



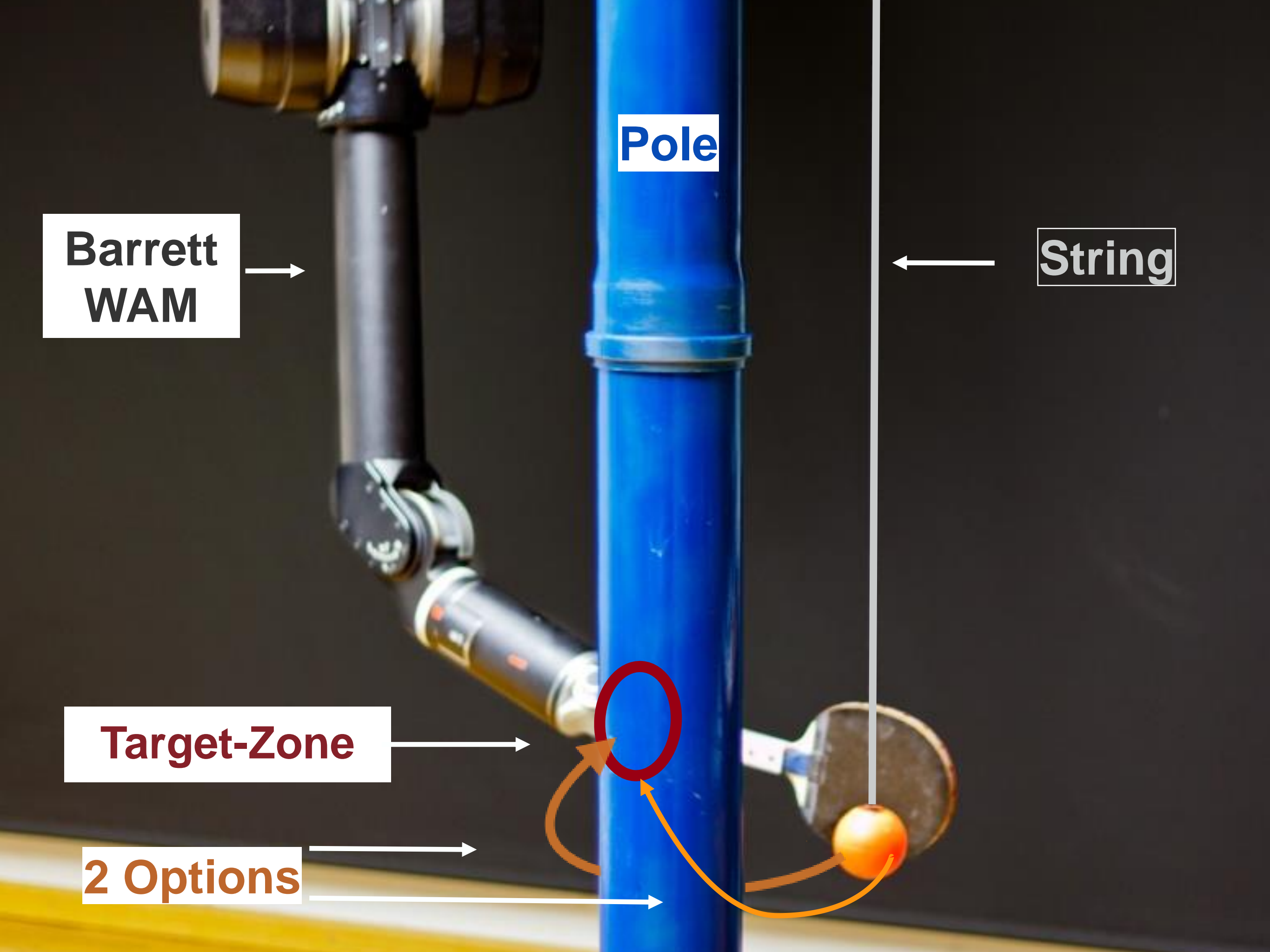
Pole

**Barrett
WAM**

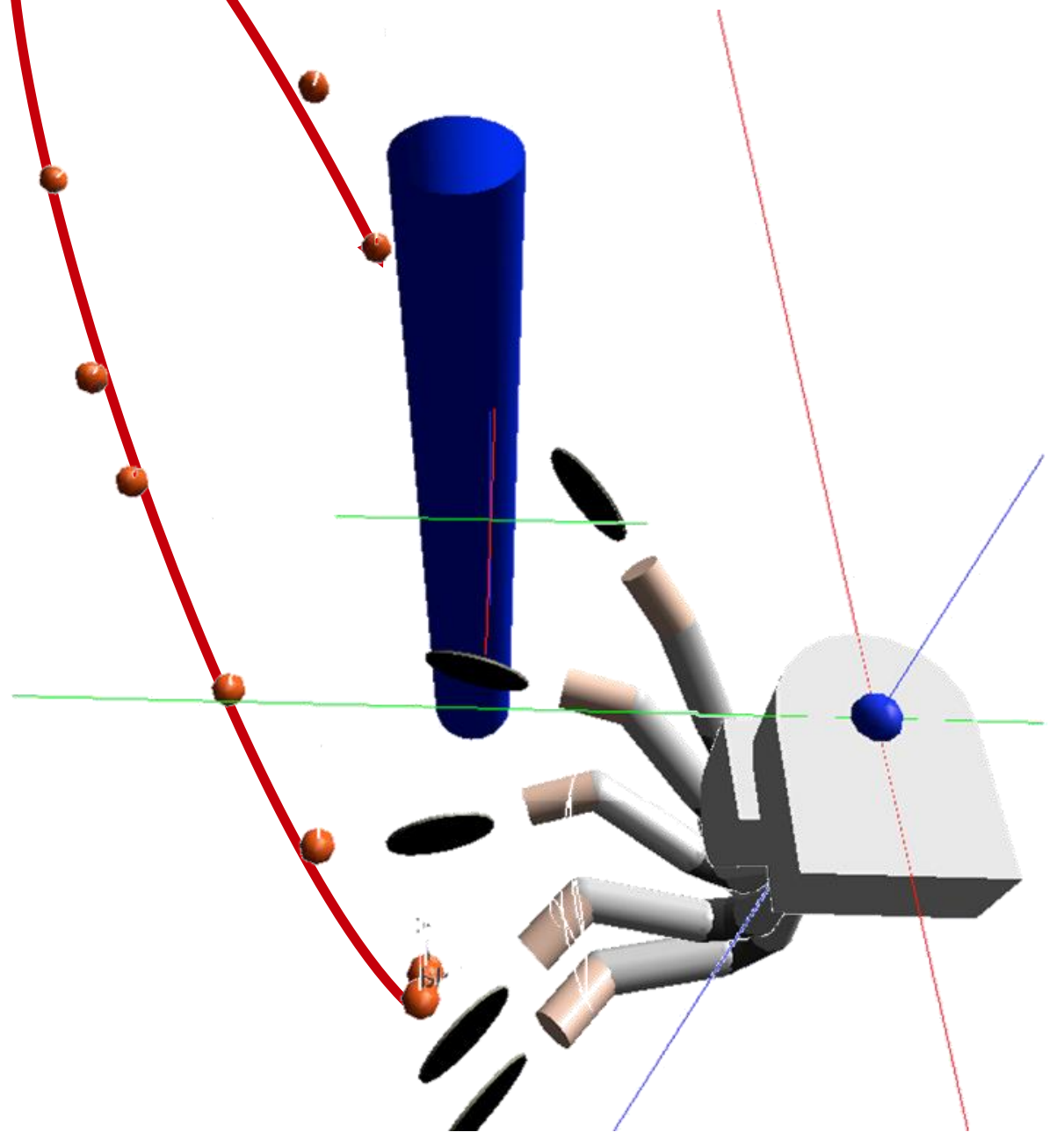
String

Target-Zone

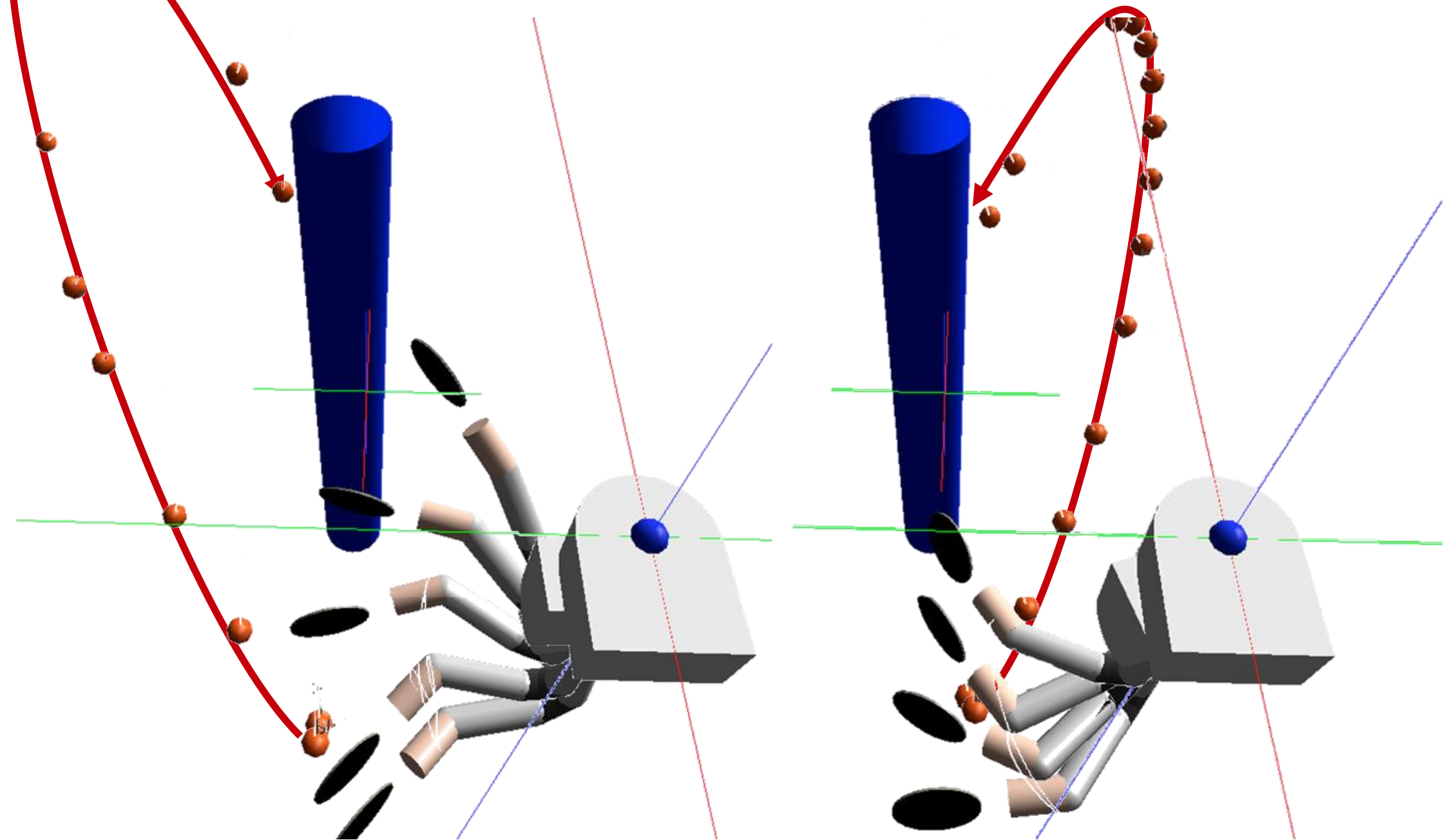
2 Options



Tetherball



Tetherball



HiREPS learns distinct solutions.

Conclusion

Probabilistic Policy Search Methods

- Policy update reduces to **weighted maximum likelihood estimates** of the parameters
- Any type of **structured policy** can be used (e.g. mixture model)
- Weights are specified by **exponential transformation** of the returns
- REPS **optimizes the temperature** of this transformation to match a desired Kullback-Leibler divergence
- **Contextual policy search** is a powerful tool for multi-task learning