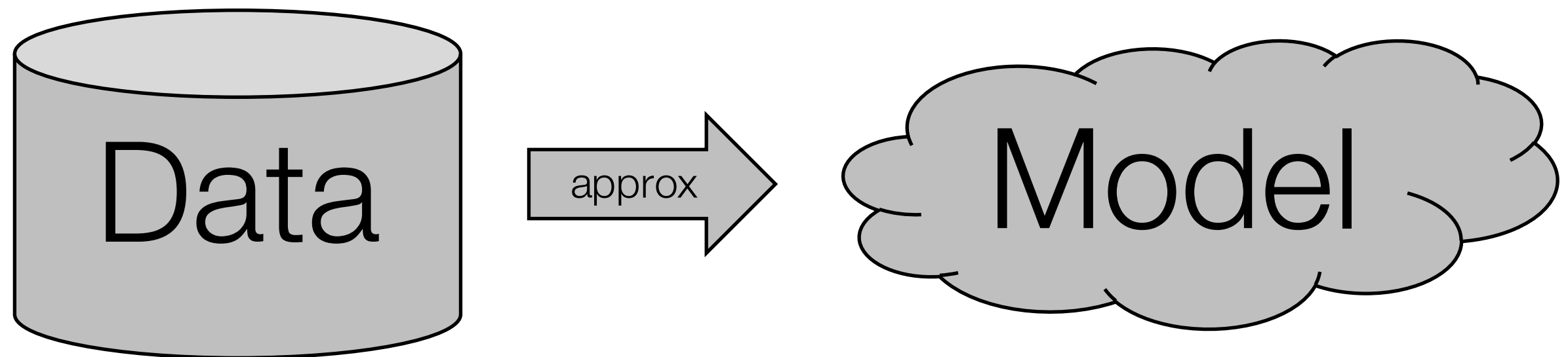




Reinforcement Learning: A Tutorial

Jan Peters

Supervised/Unsupervised Learning



Motivation for optimal decision making in robotics



Typically, **supervised learning is not enough**

Imperfect demonstrations

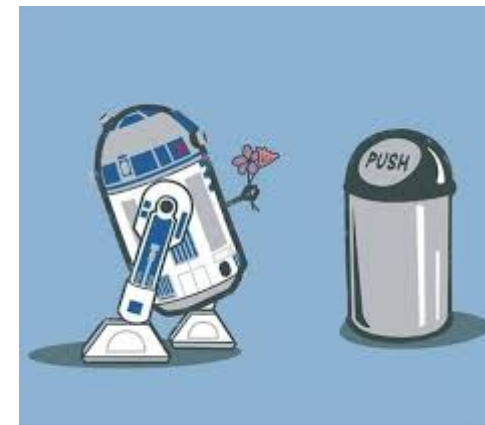
We cannot demonstrate everything!

ML systems need **self-improvement!**

The system explores by trial and error

We give evaluative feedback  reward

Today, we are going to look at the problem of how to **derive optimal actions that maximize long-term reward**



Exploration



Reward

Note:

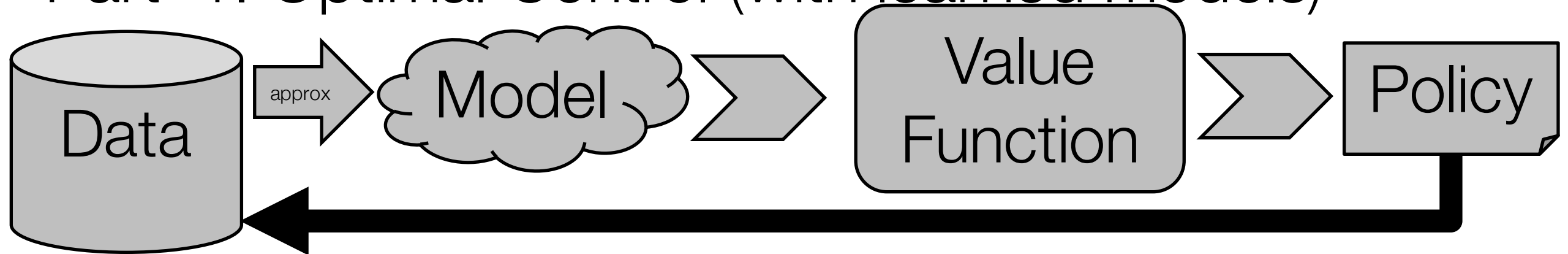
reward = - cost

Max(reward) = Min(cost)

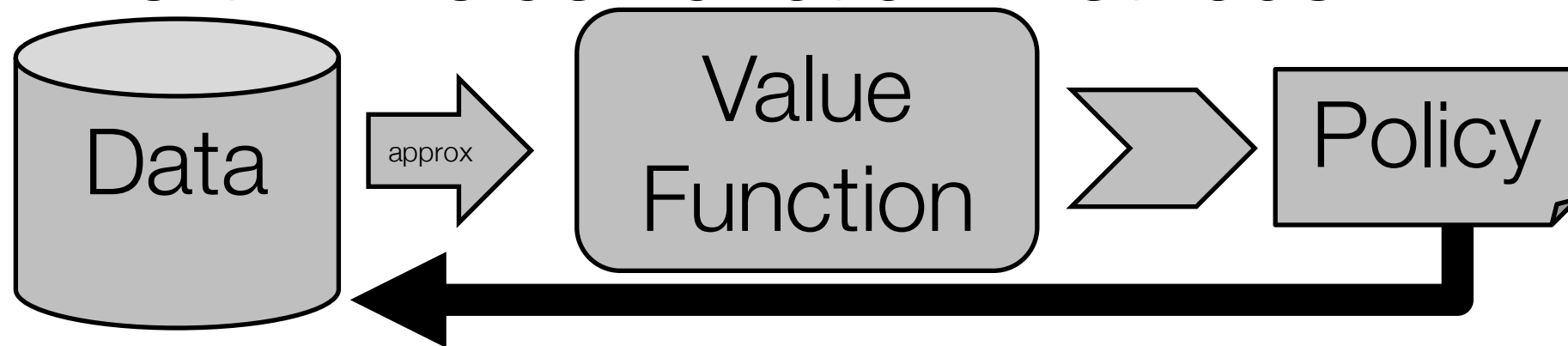
Reinforcement Learning



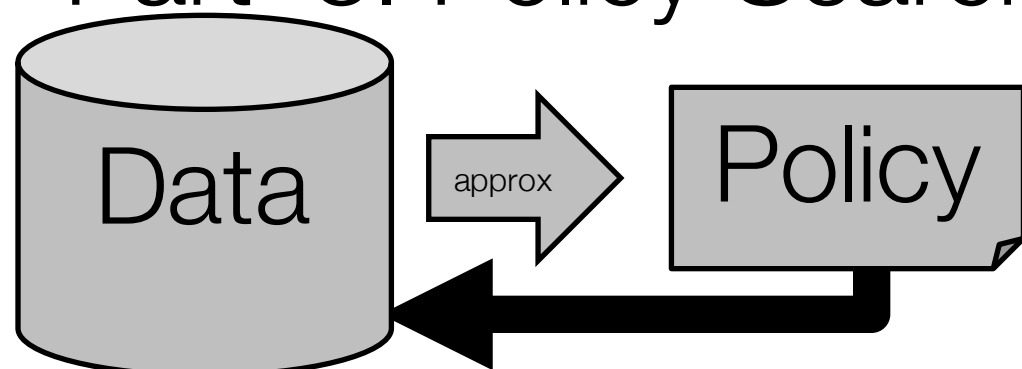
Part 1. Optimal Control (with learned models)



Part 2. Value Function Methods



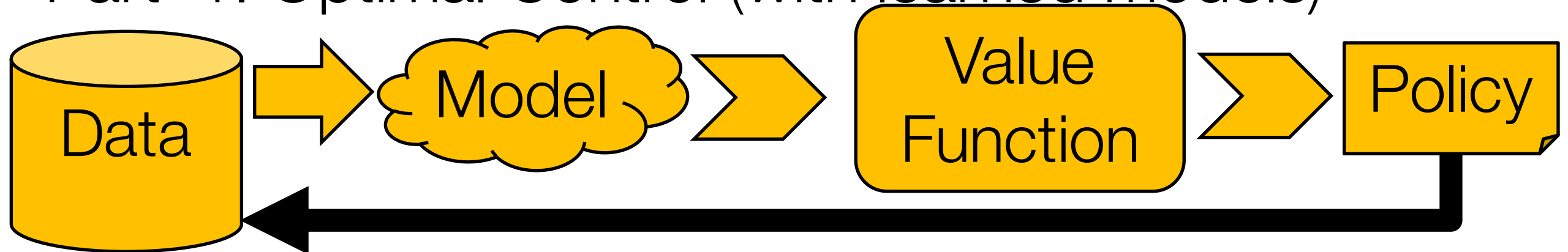
Part 3. Policy Search



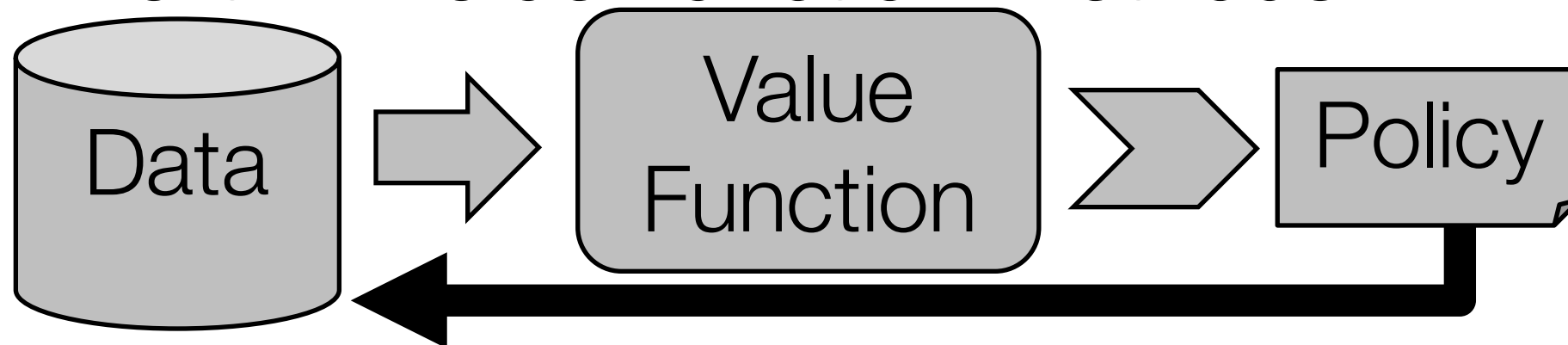
Reinforcement Learning



Part 1. Optimal Control (with learned models)



Part 2. Value Function Methods



Part 3. Policy Search

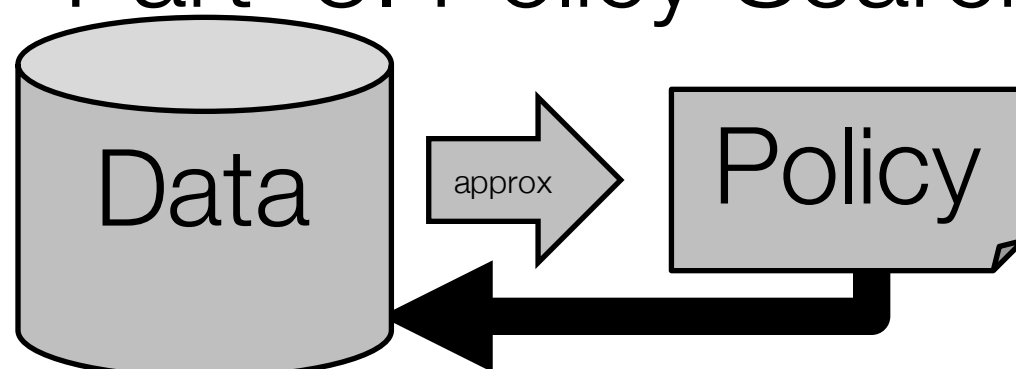
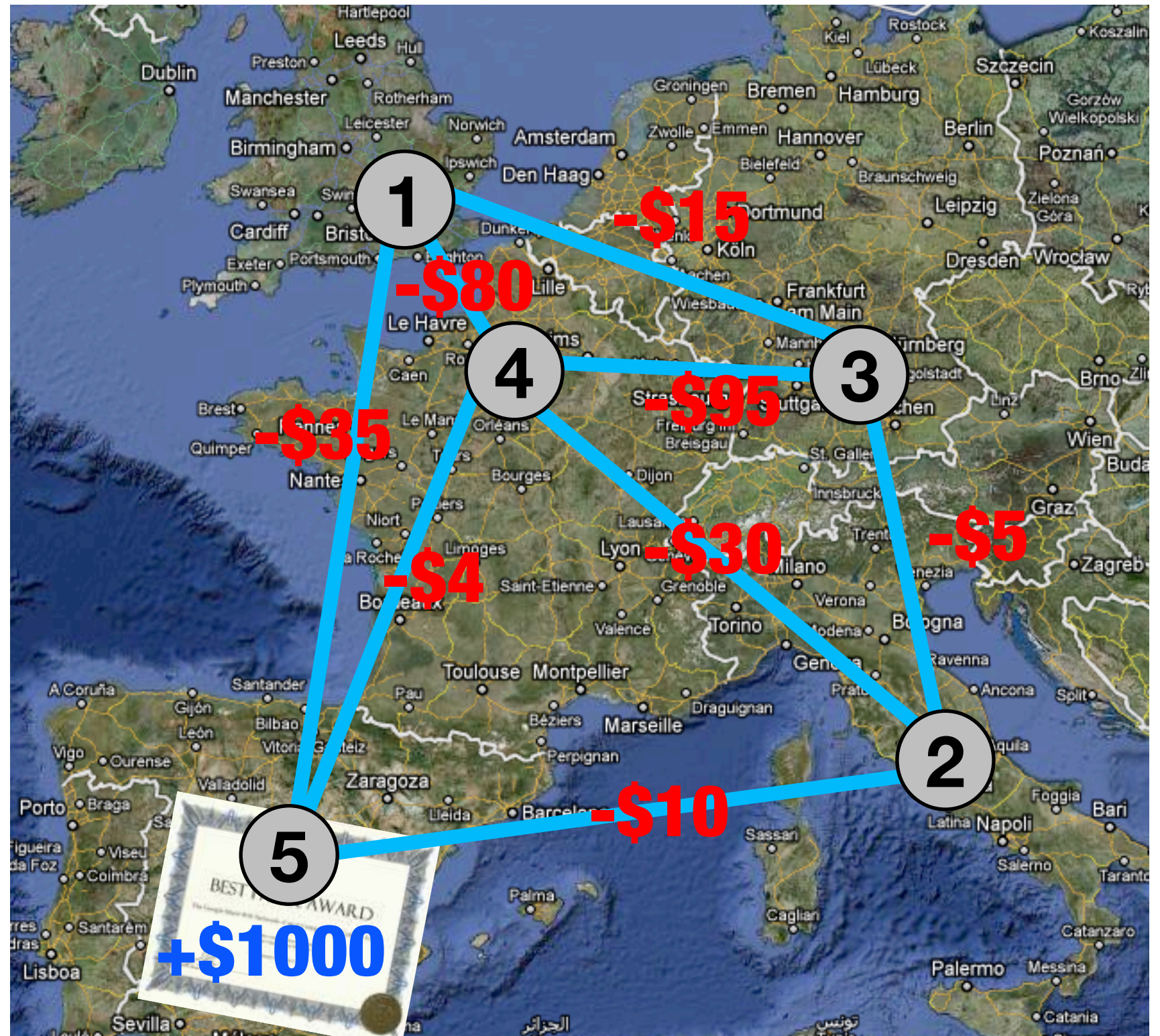


Illustration of basic idea...



**You have won
a Best-Paper
Award in
Madrid!**

**What is the
Optimal
Policy to
Collect it?**



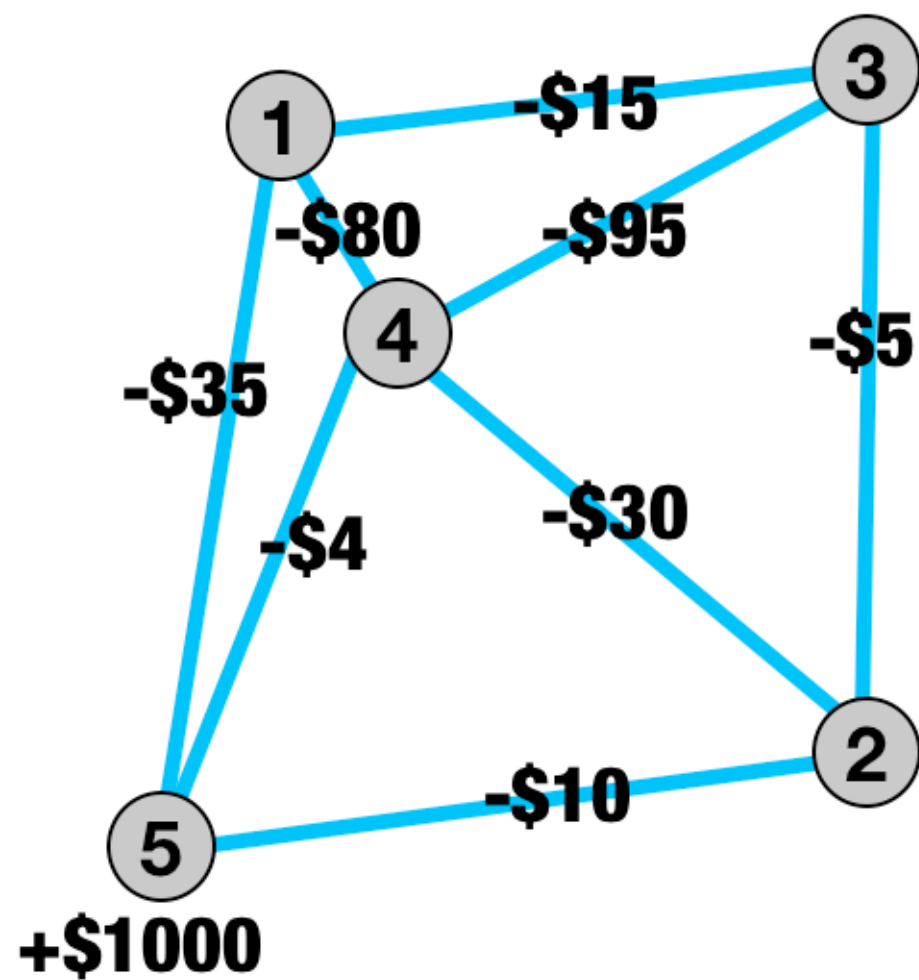
Dynamic Programming



*“An optimal sequence of controls in a multistage optimization problem has the property that **whatever the initial stage, state and controls are**, the **remaining controls** must constitute **an optimal sequence of decisions for the remaining problem** with stage and state resulting from previous controls considered as initial conditions.”*



Let's Try this Example!



	T-4	T-3	T-2	T-1	T	
					0	1
					0	2
					0	3
					0	4
					1000	5

Arrows indicating transitions between states (rows) over time (columns):

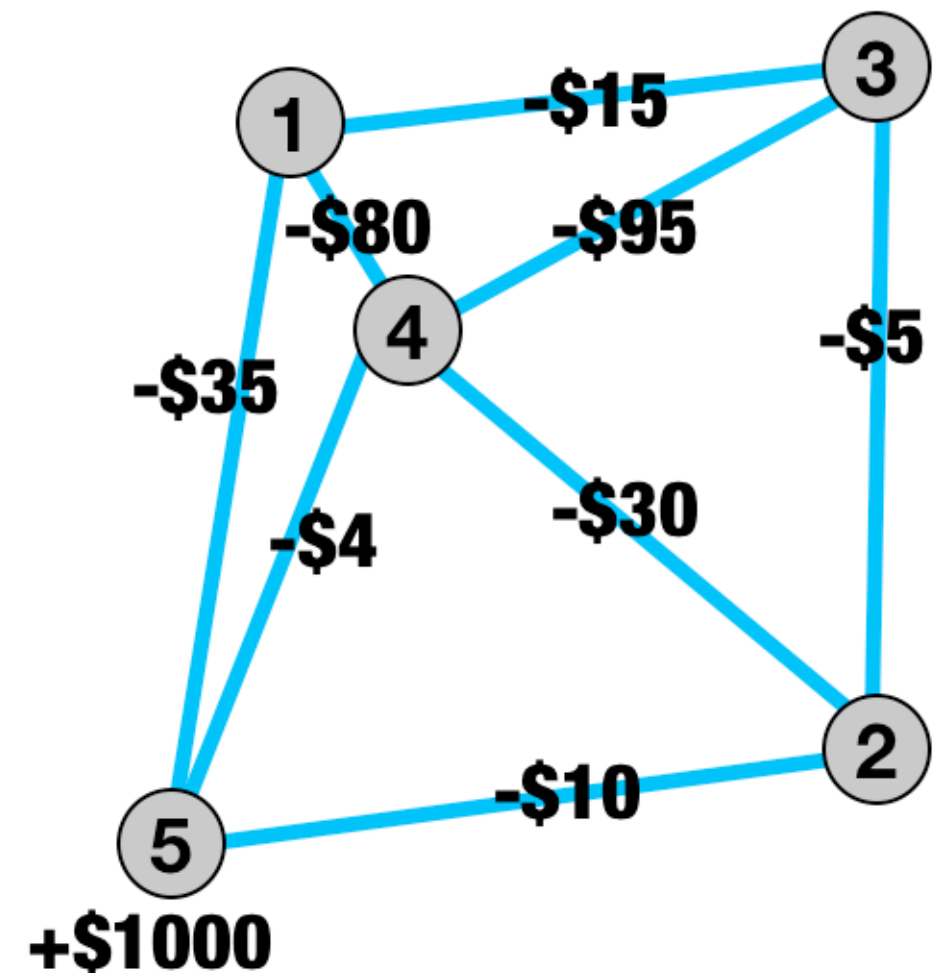
- Green arrows: 1 → 2, 2 → 3, 3 → 4, 4 → 5
- Red arrows: 1 → 3, 2 → 4, 3 → 5
- Blue arrows: 1 → 4, 2 → 5
- Purple arrows: 1 → 5, 2 → 5
- Yellow arrows: 1 → 5, 2 → 5

Markov Decision Problems (MDP)



A stationary **MDP** is defined by:

- its state space $s \in \mathcal{S}$
- its action space $a \in \mathcal{A}$
- its transition dynamics $\mathcal{P}(s_{t+1} | s_t, a_t)$
- its reward function $r(s, a)$
- and its initial state probabilities $\mu_0(s)$



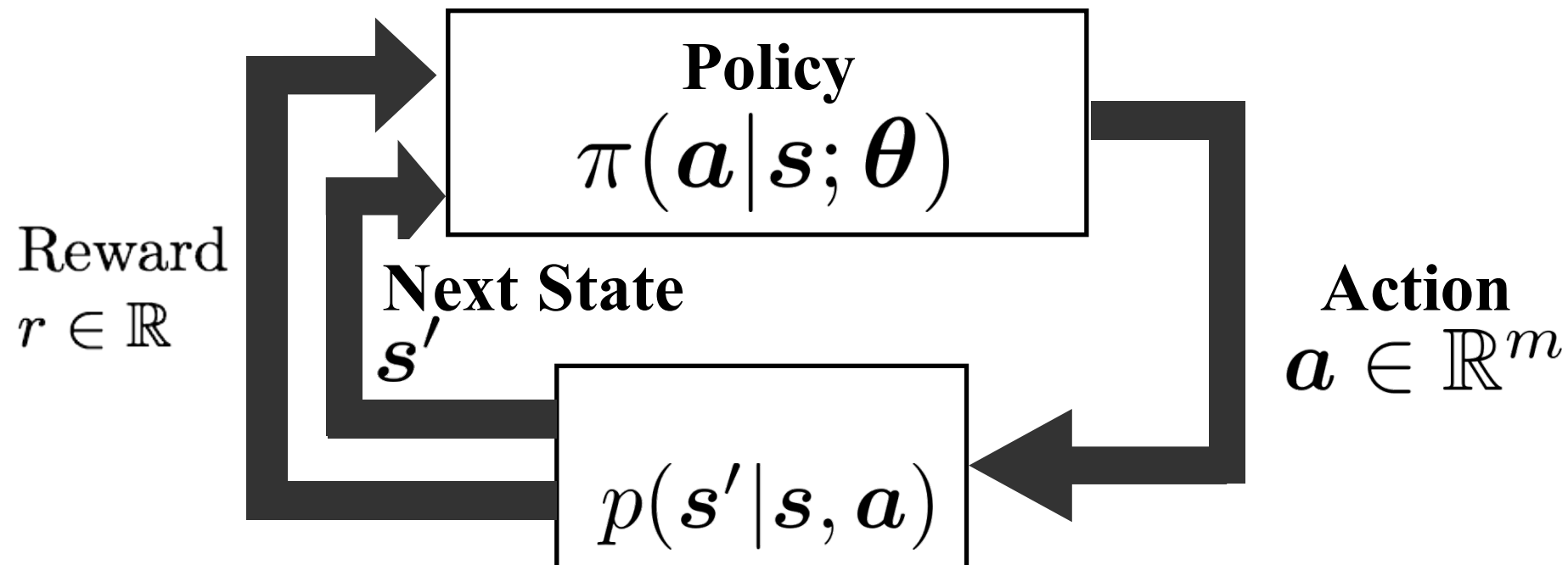
Markov property:

$$\mathcal{P}(s_{t+1} | s_t, a_t, s_{t-1}, a_{t-1}, \dots) = \mathcal{P}(s_{t+1} | s_t, a_t)$$

- Transition dynamics depends on only on the current time step



Basic Reinforcement Learning Loop:



Goal: Maximize the expected long-term reward

$$J_{\theta} = \mathbb{E}_{\mu_0, \mathcal{P}, \pi} \left[\sum_{t=1}^{T-1} \gamma^t r(s_t, a_t) \right]$$

10 discount factor $0 \leq \gamma \leq 1$



Algorithmic Description of Value Iteration

Init: $V_T^*(s) \leftarrow r_T(s)$, $t = T$

Repeat $t = t - 1$

Compute Q-Function for time step t (for each state action pair)

$$Q_t^*(s, a) = r_t(s, a) + \gamma \sum_{s'} P_t(s'|s, a) V_{t+1}^*(s')$$

Compute V-Function for time step t (for each state)

$$V_t^*(s) = \max_a Q_t^*(s, a)$$

Until $t = 1$

Return: Optimal policy for **each time step**

$$\pi_t^*(s) = \operatorname{argmax}_a Q_t^*(s, a)$$



What if you

„Bellman Equation“ (Bellman Principle of Optimality)

$$V^*(\mathbf{s}) = \max_{\mathbf{a}} \left(r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathcal{P}} [V^*(\mathbf{s}') | \mathbf{s}, \mathbf{a}] \right)$$

➡ Iterating the Bellman Equation converges to the stationary value function V^*

Alternatively, we can write this in Q-Functions

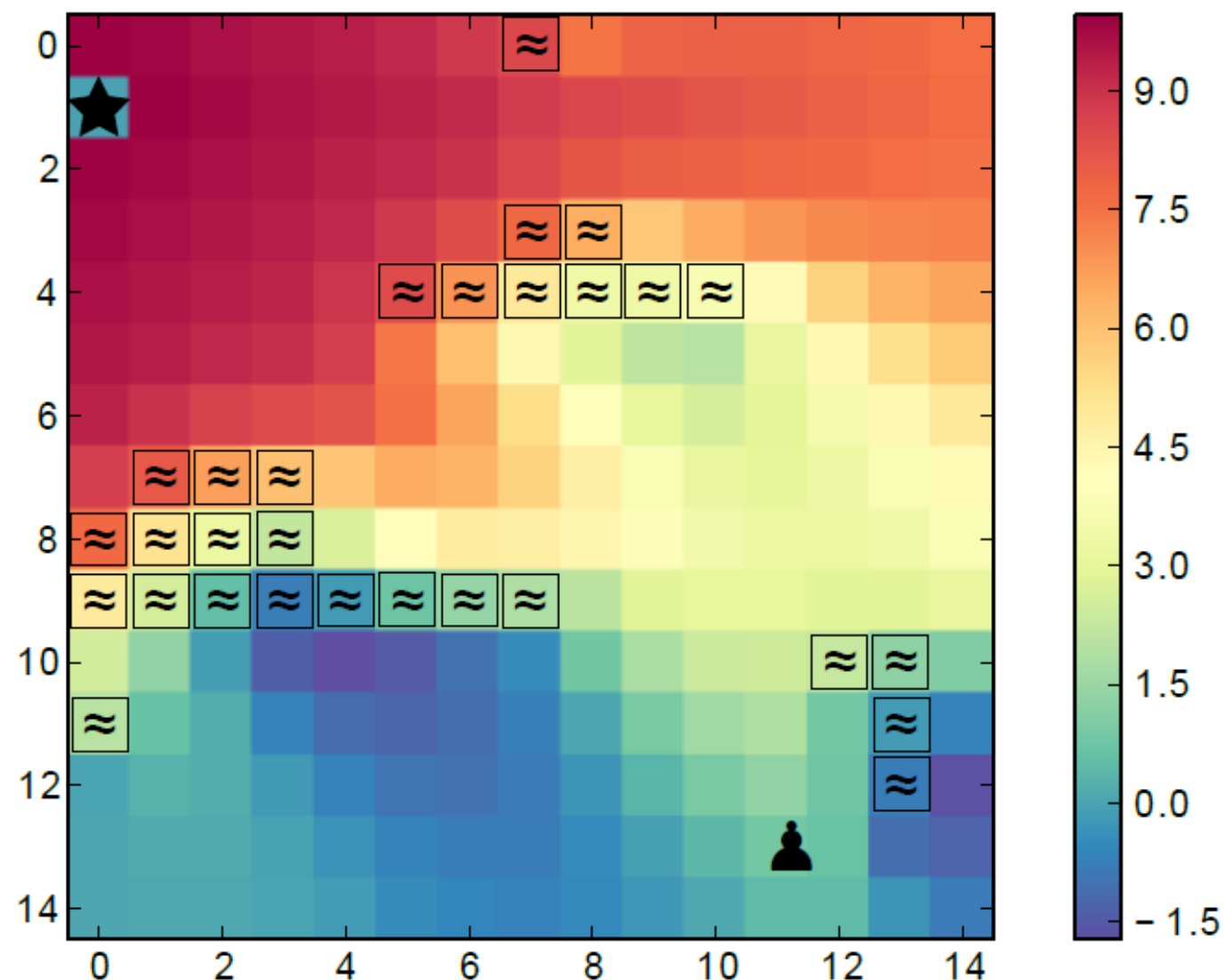
$$Q^*(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathcal{P}} [\max_{\mathbf{a}'} Q^*(\mathbf{s}', \mathbf{a}') | \mathbf{s}, \mathbf{a}]$$

If your life is infinite: Stationary value functions



An Illustration...

Policy always goes directly to the star
Going through puddles is punished





What if the max is expensive?

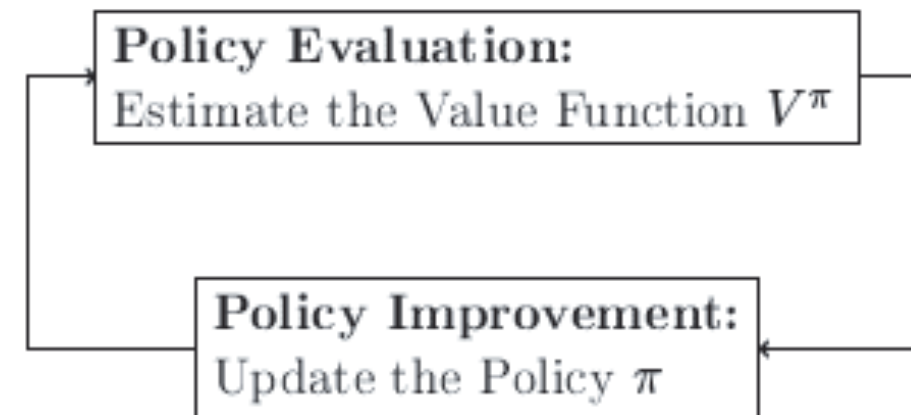
Typically done iteratively:

- **Policy Evaluation:**

Estimate quality of states (and actions) with current policy

- **Policy Improvement:**

Improve policy by taking actions with the highest quality



Such iterations are called **Policy Iteration**.

A Special MDP: Linear Quadratic Gaussian Systems



An **LQR** system is defined as

- its state space $\mathbf{x} \in \mathbb{R}^n$ (note: same as \mathcal{S})
- its action space $\mathbf{u} \in \mathbb{R}^m$ (note: same as \mathcal{U})
- its (possibly time-dependent) **linear transition dynamics with Gaussian noise**

$$p_t(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t + \mathbf{b}_t, \Sigma_t)$$

- its **quadratic** reward function

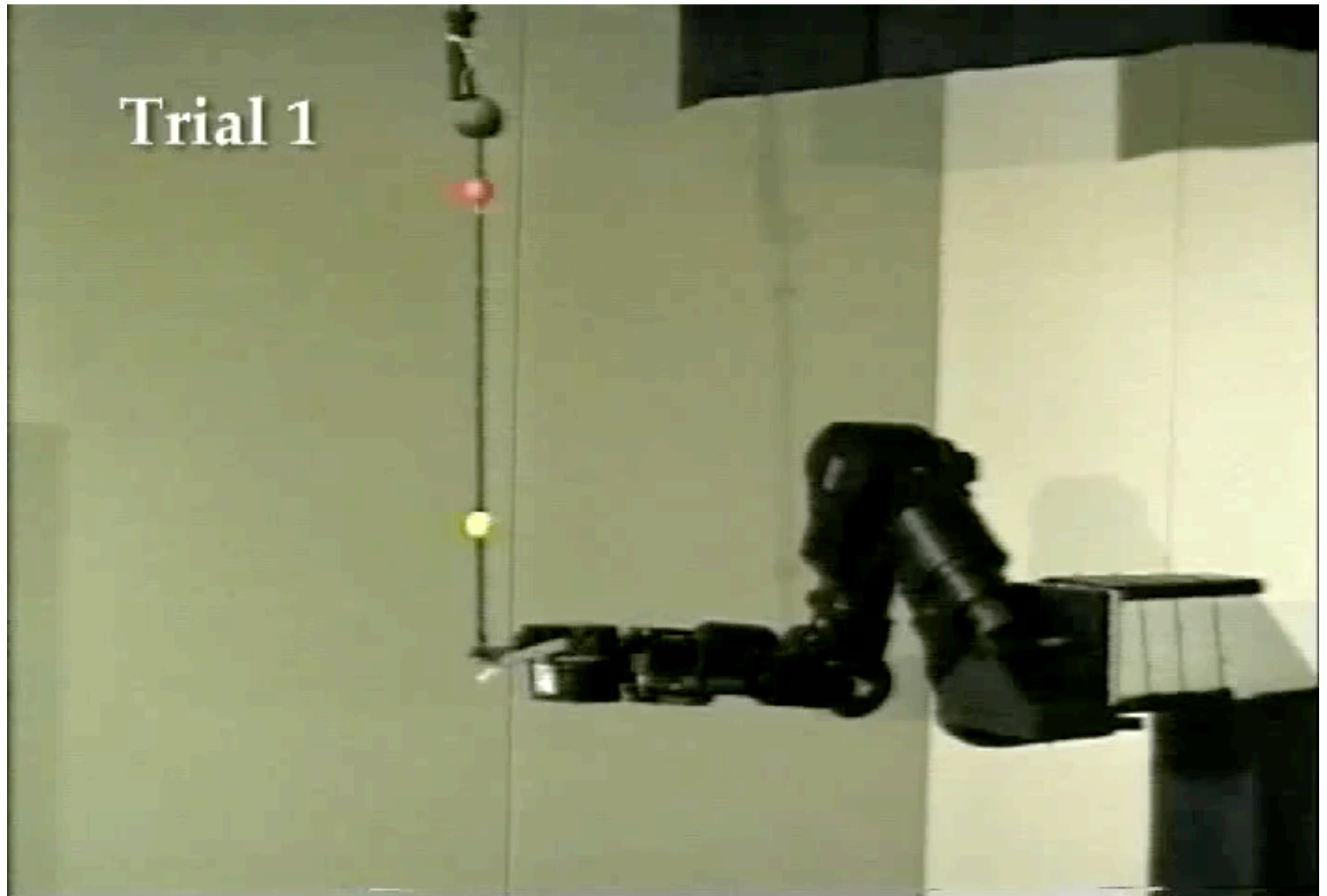
$$r_t(\mathbf{x}, \mathbf{u}) = (\mathbf{x} - \mathbf{r}_t)^T \mathbf{R}_t (\mathbf{x} - \mathbf{r}_t) + \mathbf{u}_t^T \mathbf{H}_t \mathbf{u}_t$$

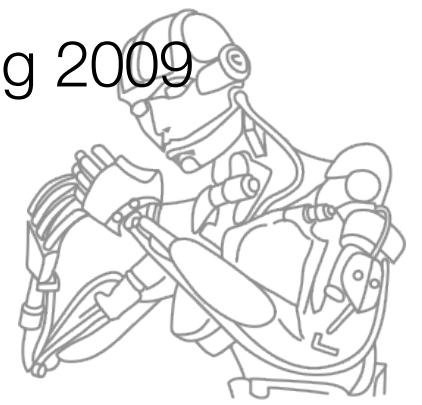
$$r_T(\mathbf{x}) = (\mathbf{x} - \mathbf{r}_T)^T \mathbf{R}_T (\mathbf{x} - \mathbf{r}_T)$$

- and its initial state density

$$\mu_0(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mu_0, \Sigma_0)$$

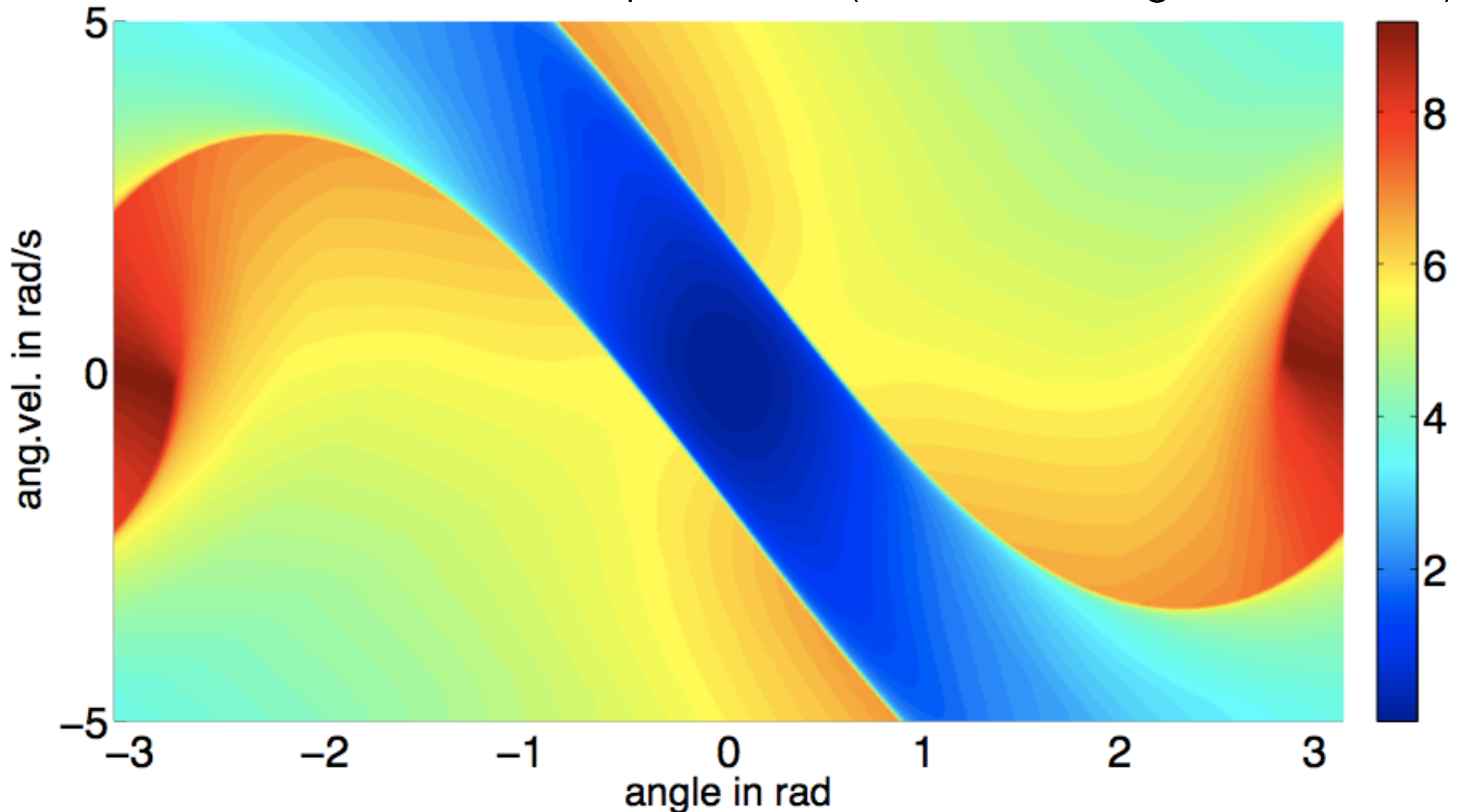
A Special MDP: Linear Quadratic Gaussian Systems





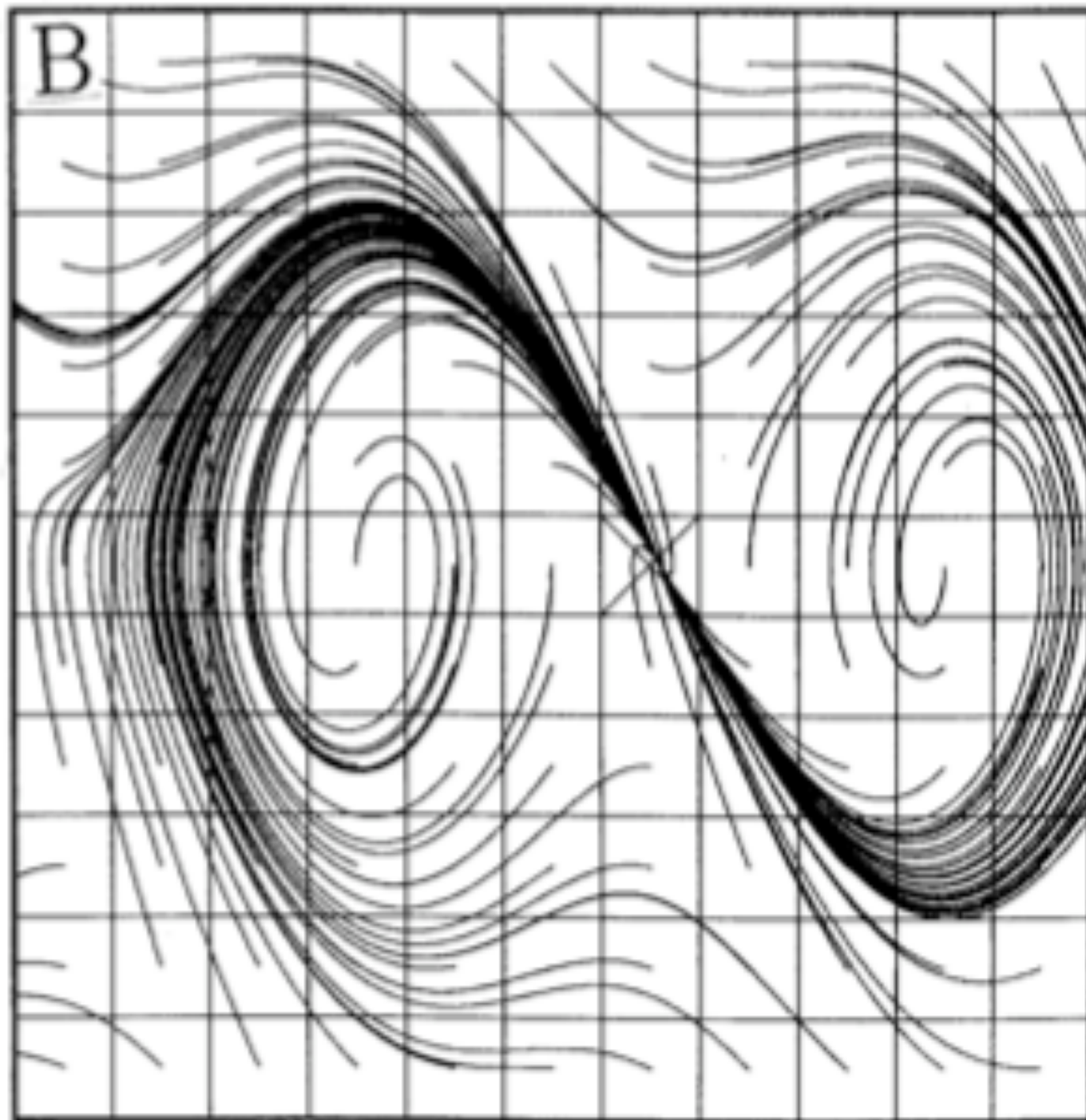
What's wrong with LQR?

Value function for the inverted pendulum (on costs = negative rewards)



Highly non-linear function (certainly not quadratic)

Possible: Learn Solutions only where needed!



**If you know places
where we start...**

**... we can just look
ahead and
approximate the
solution locally
around an initial
trajectory**



Local Solutions by Linearizations

Every smooth function can be modeled with a Taylor expansion

$$f(\mathbf{x}) = f(\mathbf{a}) + \left. \frac{df}{d\mathbf{x}} \right|_{\mathbf{x}=\mathbf{a}} (\mathbf{x} - \mathbf{a}) + \frac{1}{2} (\mathbf{x} - \mathbf{a})^T \left. \frac{d^2 f}{d\mathbf{x}^2} \right|_{\mathbf{x}=\mathbf{a}} (\mathbf{x} - \mathbf{a}) + \dots$$

Hence, we can also **approximate the (learned) forward dynamics by linearizing** at the point $(\tilde{\mathbf{x}}_t, \tilde{\mathbf{u}}_t)$

$$\begin{aligned} \mathbf{x}_{t+1} = f_t(\mathbf{x}_t, \mathbf{u}_t) &\approx f(\tilde{\mathbf{x}}_t, \tilde{\mathbf{u}}_t) + \frac{df}{ds} (\mathbf{x}_t - \tilde{\mathbf{x}}_t) + \frac{df}{du} (\mathbf{u}_t - \tilde{\mathbf{u}}_t) \\ &= \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t + \mathbf{b}_t \end{aligned}$$

and **approximate the (learned) reward function by a second order approximation**

$$r_t(\mathbf{s}_t, \mathbf{a}_t) \approx r(\tilde{\mathbf{x}}_t, \tilde{\mathbf{u}}_t) + \frac{dr}{d\mathbf{x}} (\mathbf{x}_t - \tilde{\mathbf{x}}_t) + (\mathbf{x}_t - \tilde{\mathbf{x}}_t)^T \frac{dr}{d\mathbf{x}d\mathbf{x}} (\mathbf{x}_t - \tilde{\mathbf{x}}_t) - \mathbf{u}_t^T \mathbf{H}_t \mathbf{u}_t$$



Local Solutions by Linearizations

So we are back to the **full linear optimal control case with...**

$$p_t(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t) = \mathcal{N}(\mathbf{x}_{t+1}|\mathbf{A}_t\mathbf{x}_t + \mathbf{B}_t\mathbf{u}_t + \mathbf{b}_t, \Sigma_t)$$

$$r(\mathbf{x}, \mathbf{u}) = -\mathbf{x}^T \mathbf{R}_t \mathbf{x} + 2\mathbf{r}_t^T \mathbf{x} - \mathbf{u}^T \mathbf{H}_t \mathbf{u} + \text{const}$$

that we know how to solve...

Hence our algorithm for **solving non-linear optimal control** is...

- 1. Backward Solution:** Compute optimal control law (i.e. Gains \mathbf{K}_t and offsets \mathbf{k}_t)
- 2. Forward Propagation:** Run simulator with optimal control law to obtain linearization points $(\tilde{\mathbf{x}}_{1:T}, \tilde{\mathbf{u}}_{1:T})$

1.If not converged, go to 1.

Application to the Swing-Up



Some interesting results (only in simulation)



Work by Emo Todorov
and Yuval Tassa
(They call basically the
same algorithm
incremental LQG, iLQG).

Note: iLQG is just a
simplification of Differential
Dynamic Programming
(Dyer & McReynolds, 1969)

Synthesis of Complex Behaviors
with
Online Trajectory Optimization

(under review)



Wrap-Up: Optimal Control

We now know how to compute **optimal policies**

Cool, that's all we need. Let's go home...

Wait, **there is a catch! Unfortunately, we can only do this in 2.5 cases**

- Discrete Systems

Easy: integrals turn into sums

...but the world is not discrete!

- Linear Systems, Quadratic Reward, Gaussian Noise (LQR)

... but the world is not linear!

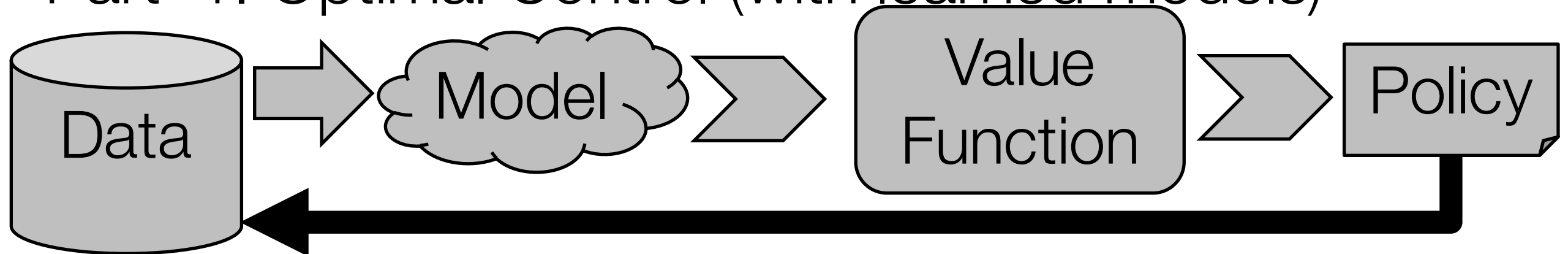
- Along an optimal trajectory – finding it is really hard!

Otherwise, we need to approximate!

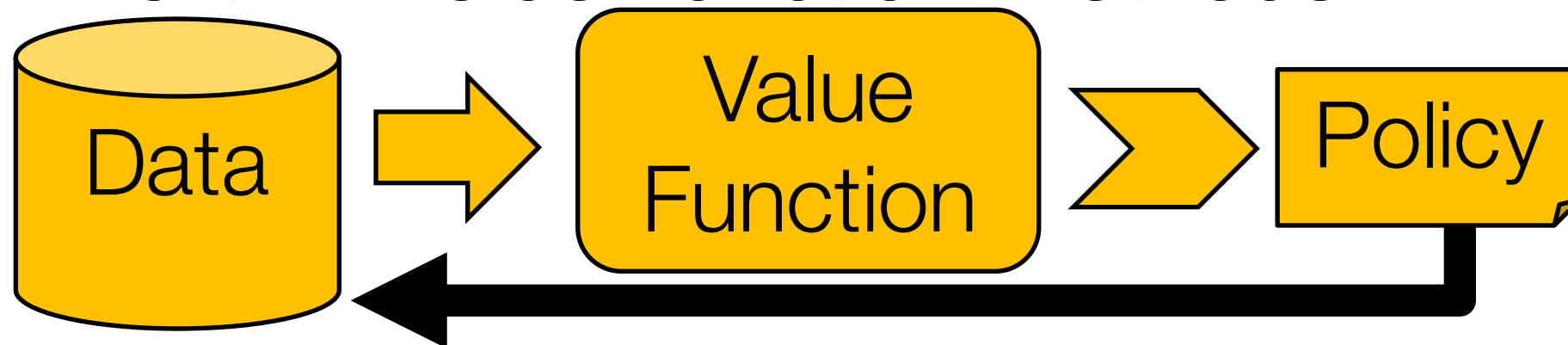
Reinforcement Learning



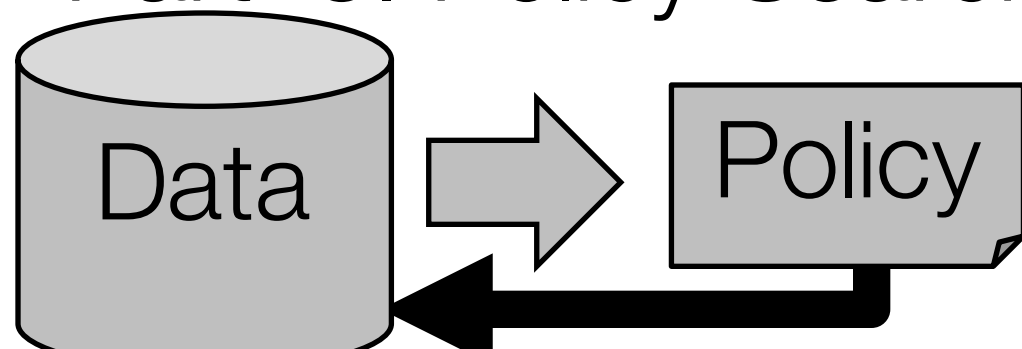
Part 1. Optimal Control (with learned models)



Part 2. Value Function Methods



Part 3. Policy Search



CHRIS ATKESON Humanoids 2016

American election:

- Clinton was model-based, and used strong predictive models of who would vote and how they would vote.
- Trump did not use any models.



Purpose of this Lecture

Often, learning a good model is too hard

- ➔ The optimization inherent in optimal control is **prone to model errors**, as the controller may achieve the objective only because model errors get exploited
- ➔ Optimal control methods based on linearization of the dynamics work only for **moderately non-linear tasks**
- ➔ (Ideally model-free) Approaches are needed that do not make any assumption on the structure of the model

Classical Reinforcement Learning:

- ➔ Solve the optimal control problem by **learning the value function, not the model!**



Markov Decision Processes (MDP)

Classical reinforcement learning is typically formulated for the infinite horizon objective

Infinite Horizon: maximize **discounted accumulated reward**

$$J_{\pi} = \mathbb{E}_{\mu_0, \mathcal{P}, \pi} \left[\sum_{t=0}^{\infty} \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right]$$

$0 \leq \gamma < 1$... discount factor (note change!)

Trades-off long term vs. immediate reward



Value functions of a policy

Value function and state-action value function of a policy can be computed iteratively

$$\begin{aligned} V^\pi(\mathbf{s}) &= \mathbb{E}_\pi \left[r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathcal{P}} [V^\pi(\mathbf{s}')] \mid \mathbf{s} \right] \\ &= \int \pi(\mathbf{a} \mid \mathbf{s}) \left(r(\mathbf{s}, \mathbf{a}) + \gamma \int \mathcal{P}(\mathbf{s}' \mid \mathbf{s}, \mathbf{a}) V^\pi(\mathbf{s}') d\mathbf{s}' \right) d\mathbf{a} \end{aligned}$$

$$\begin{aligned} Q^\pi(\mathbf{s}, \mathbf{a}) &= r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathcal{P}, \pi} \left[Q^\pi(\mathbf{s}', \mathbf{a}') \mid \mathbf{s}, \mathbf{a} \right] \\ &= r(\mathbf{s}, \mathbf{a}) + \gamma \int \mathcal{P}(\mathbf{s}' \mid \mathbf{s}, \mathbf{a}) \int \pi(\mathbf{a}' \mid \mathbf{s}') Q^\pi(\mathbf{s}', \mathbf{a}') d\mathbf{a}' d\mathbf{s}' \end{aligned}$$



Value-based Reinforcement Learning

Classical Reinforcement Learning

Updates the value function based on samples

$$\mathcal{D} = \{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i\}_{i=1\dots N}$$

We do not have a model and we do not want to learn it

Use the samples to update Q-function (or V-function)

Lets start simple:

Discrete states/actions ➡ Tabular Q-function



Temporal difference learning

Given a transition (s_t, a_t, r_t, s_{t+1}) , we want to update the V-function

- Use the estimate of the current value: $V(s_t)$
- 1-step prediction of the current value: $\hat{V}(s_t) = r_t + \gamma V(s_{t+1})$
- 1-step prediction error (called temporal difference (TD) error)

$$\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$$

Update current value with the temporal difference error

$$V_{\text{new}}(s_t) = V(s_t) + \alpha \delta_t = (1 - \alpha)V(s_t) + \alpha(r_t + \gamma V(s_{t+1}))$$



Temporal difference learning

The **TD error**

$$\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$$

compares the **one-time step lookahead prediction**

$$\hat{V}(s_t) = r_t + \gamma V(s_{t+1})$$

with the **current estimate** of the value function $V(s_t)$

⇒ if $\hat{V}(s_t) > V(s_t)$ than $V(s_t)$ is increased

⇒ if $\hat{V}(s_t) < V(s_t)$ than $V(s_t)$ is decreased



Algorithmic Description of TD Learning

Init: $V_0^*(s) \leftarrow 0$

Repeat $t = t + 1$

Observe transition (s_t, a_t, r_t, s_{t+1})

Compute TD error $\delta_t = r_t + \gamma V_t(s_{t+1}) - V_t(s_t)$

Update V-Function $V_{t+1}(s_t) = V_t(s_t) + \alpha \delta_t$

until convergence of V

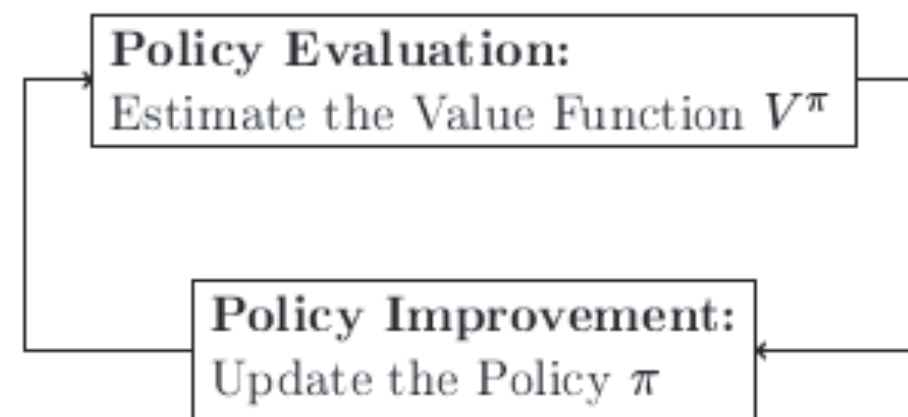
➡ Used to compute Value function of behavior policy

➡ Sample-based version of policy evaluation



Temporal difference learning for control

So far: Policy evaluation with TD methods



Can we also do the policy improvement step **with samples**?

Yes, but we need to enforce exploration!

Epsilon-Greedy Policy: $\pi(a|s) = \begin{cases} 1 - \epsilon + \epsilon/|\mathcal{A}|, & \text{if } a = \operatorname{argmax}_{a'} Q^\pi(s, a') \\ \epsilon/|\mathcal{A}|, & \text{otherwise} \end{cases}$

Soft-Max Policy: $\pi(a|s) = \frac{\exp(\beta Q(s, a))}{\sum_{a'} \exp(\beta Q(s, a'))}$

33 ➡ **Do not always take greedy action**



Temporal difference learning for control

Update equations for **learning the Q-function** $Q(s, a)$

$$Q_{t+1}(s_t, a_t) = Q_t(s_t, a_t) + \alpha \delta_t, \quad \delta_t = r_t + \gamma Q_t(s_{t+1}, a?) - Q_t(s_t, a_t)$$

Two different methods to estimate $a?$

Q-learning: $a? = \operatorname{argmax}_a Q_t(s_{t+1}, a)$

Estimates Q-function of optimal policy

Off-policy samples: $a? \neq a_{t+1}$

SARSA: $a? = a_{t+1}$, where $a_{t+1} \sim \pi(a|s_{t+1})$

Estimates Q-function of exploration policy

On-policy samples

Note: The policy for generating the actions depends on the Q-function → non-stationary policy



Approximating the Value Function

In the continuous case, we need to approximate the V-function (except for LQR)

Lets keep it simple, we use **a linear model** to represent the V-function

$$V^{\pi}(s) \approx V_{\omega}(s) = \phi^T(s)\omega$$

How can we find the parameters ω ?

➡ Again with **Temporal Difference Learning**



TD-learning with Function Approximation

Derivation:

Use the **recursive definition of V-function**:

$$\text{MSE}(\omega) \approx \text{MSE}_{\text{BS}}(\omega) = 1/N \sum_{i=1}^N \left(\hat{V}^{\pi}(\mathbf{s}_i) - V_{\omega}(\mathbf{s}_i) \right)^2$$

with $\hat{V}^{\pi}(\mathbf{s}) = \mathbb{E}_{\pi} \left[r(\mathbf{s}, \mathbf{a}) + \mathbb{E}_{\mathcal{P}} [V_{\omega_{\text{old}}}(\mathbf{s}') | \mathbf{s}, \mathbf{a}] \right]$

➡ **Bootstrapping (BS)**: Use the old approximation to get the target values for a new approximation

How can we **minimize** this function ?

Lets use **stochastic gradient descent**



Temporal difference learning

Stochastic gradient descent on our error function MSE_{BS}

$$\begin{aligned} MSE_{BS,t}(\omega) &= 1/N \sum_{i=1}^N \left(\hat{V}(s_t) - V_{\omega}(s_i) \right)^2 \\ &= 1/N \sum_{i=1}^N \left(r_i + \gamma V_{\omega_t}(s'_i) - V_{\omega}(s_i) \right)^2 \end{aligned}$$

Update rule (for current time step $t, V_{\omega}(s) = \phi^T(s)\omega$)

$$\begin{aligned} \omega_{t+1} &= \omega_t + \alpha_t \left. \frac{dMSE_{BS}}{d\omega} \right|_{\omega=\omega_t} \\ \omega_{t+1} &= \omega_t + \alpha \left(r(s_t, a_t) + \gamma V_{\omega_t}(s_{t+1}) - V_{\omega_t}(s_t) \right) \phi^T(s_t) \\ &= \omega_t + \alpha \delta_t \phi^T(s_t) \end{aligned}$$

with $\delta_t = r(s_t, a_t) + \gamma V_{\omega_t}(s_{t+1}) - V_{\omega_t}(s_t)$



Temporal difference learning

TD with function approximation

$$\omega_t = \omega_t + \alpha \delta_t \phi^T(s_t)$$

Difference to discrete algorithm:

- ➡ TD-error is correlated with the feature vector
- ➡ Equivalent if tabular feature coding is used, i.e., $\phi(s_i) = e_i$

Similar update rules can be obtained for SARSA and Q-learning

$$\omega_{t+1} = \omega_t + \alpha \left(r(s_t, a_t) + \gamma Q_{\omega_t}(s_{t+1}, a?) - Q_{\omega_t}(s_t, a_t) \right) \phi^T(s_t, a_t)$$

$$\text{where } Q_{\omega}(s, a) \approx \phi^T(s, a) \omega$$



Temporal difference learning

Some remarks on temporal difference learning:

- ➔ Its **not a proper** stochastic gradient descent!!
- ➔ **Why?** Target values $\hat{V}^\pi(s)$ **change after each parameter update!**

We ignore the fact that $\hat{V}^\pi(s)$ also depends on ω
- ➔ **Side note:** This „ignorance“ actually introduces a bias in our optimization, such that we are optimizing a different objective than the *MSE*
- ➔ In certain cases, we also get **divergence** (e.g. off-policy samples)
- ➔ TD-learning is very fast in terms of computation time $O(\text{\#features})$, but not data-efficient ➔ **each sample is just used once!**

Successful examples



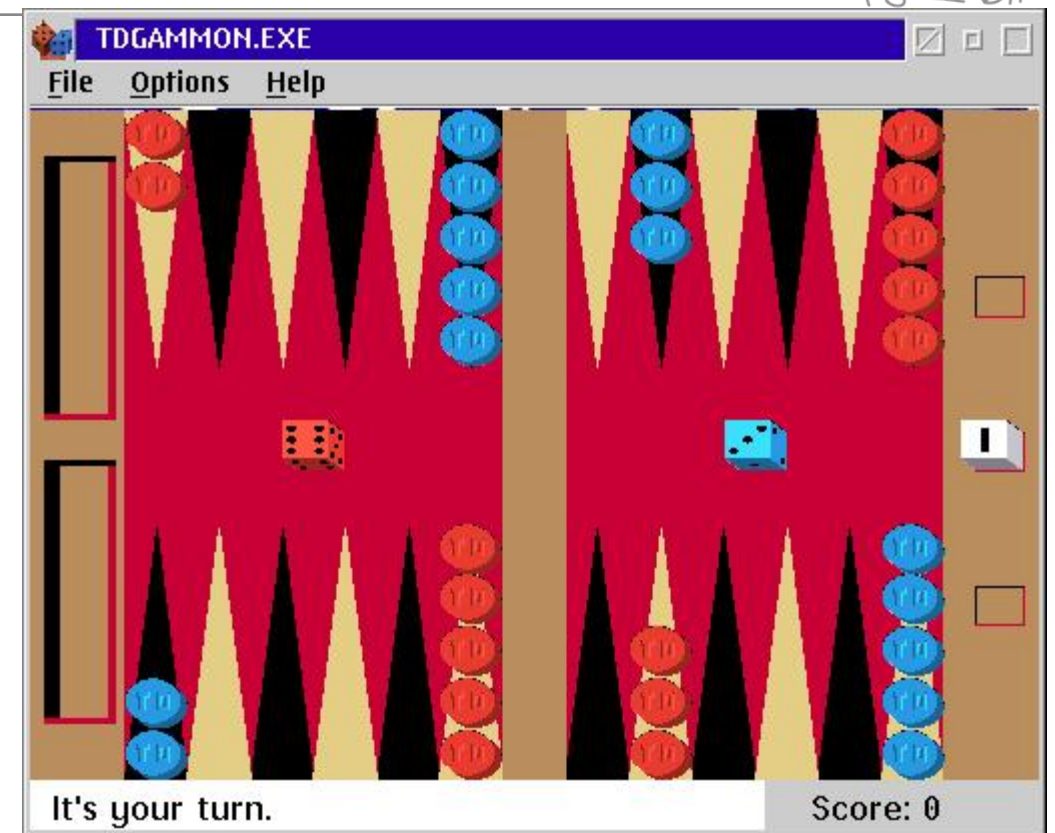
Linear function approximation

Tetris, Go

Non-linear function approximation

TD Gammon (Worldchampion level)

Atari Games (learning from raw pixel input)





Batch-Mode Reinforcement Learning

Online methods are typically **data-inefficient** as they use each data point only once

$$D = \left\{ \mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i \right\}_{i=1 \dots N}$$

Can we **re-use the whole „batch“** of data to **increase data-efficiency**?

- **Least-Squares Temporal Difference (LSTD) Learning**
- **Fitted Q-Iteration**

➡ Computationally **much more expensive** than TD-learning!



Fitted Q-iteration

In Batch-Mode RL it is also much easier to use **non-linear function approximators**

- Many of them **only exists in the batch setup**, e.g. regression trees
- **No catastrophic forgetting**, e.g., for neural networks.
- Strong divergence problems, fixed for Neural Networks by ensuring that there is a goal state where the Q-Function value is always zero (see Lange et al. below).

Fitted Q-iteration uses non-linear function approximators for **approximate value iteration**.

Ernst, Geurts and Wehenkel, *Tree-Based Batch Mode Reinforcement Learning*, JMLR 2005

Lange, Gabel and Riedmiller. *Batch Reinforcement Learning*, *Reinforcement Learning: State of the Art*



Fitted Q-iteration

Given: Dataset $D = \left\{ \mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i \right\}_{i=1 \dots N}$

Algorithm:

Initialize $Q^{[0]}(\mathbf{s}, \mathbf{a}) = 0$, input data: $\mathbf{X} = \begin{bmatrix} \mathbf{s}_1^T & \mathbf{a}_1^T \\ \vdots & \\ \mathbf{s}_N^T & \mathbf{a}_N^T \end{bmatrix}$

for $k = 1$ to L

Generate target values: $\tilde{q}_i^{[k]} = r_i + \gamma \max_{\mathbf{a}'} Q^{[k-1]}(\mathbf{s}'_i, \mathbf{a}')$

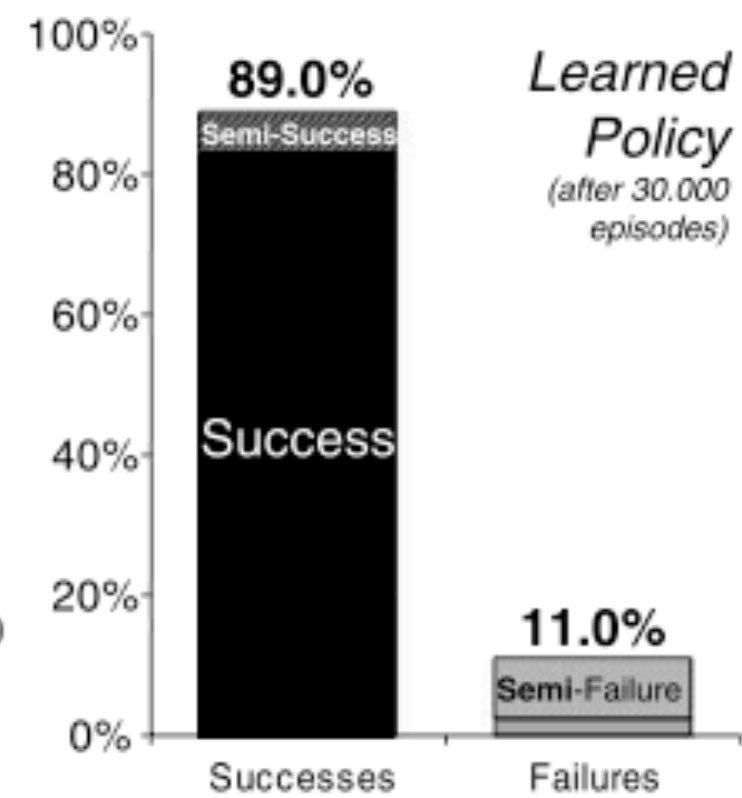
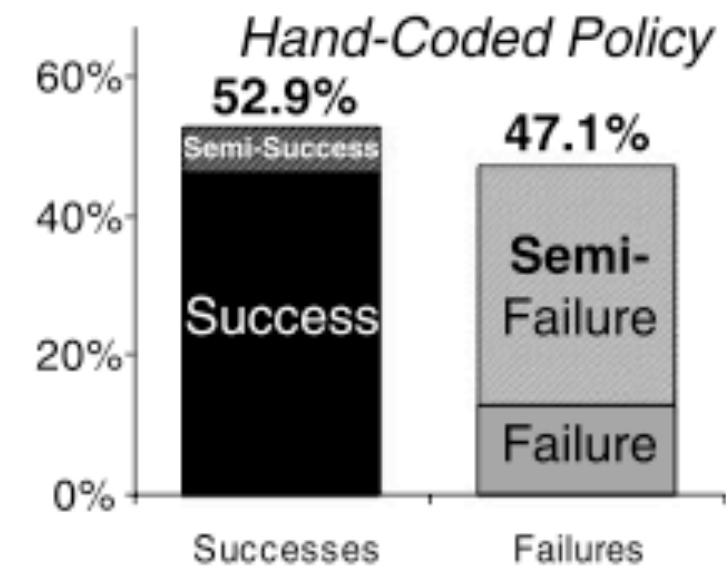
Learn new Q-function: $Q^{[k]}(\mathbf{s}, \mathbf{a}) \leftarrow \text{Regress}(\mathbf{X}, \tilde{\mathbf{q}}^{[k]})$

end

➔ Like Value-Iteration, but we use supervised learning methods to approximate the Q-function at each iteration k



Learning Robot Soccer





Value Function Methods

- ➔ ... have been the driving reinforcement learning approach in the 1990s.
- ➔ You can do loads of cool things with them: Learn Chess at professional level, learn **Backgammon and Checkers at Grandmaster-Level** ... and winning the **Robot Soccer Cup** with a minimum of man power.

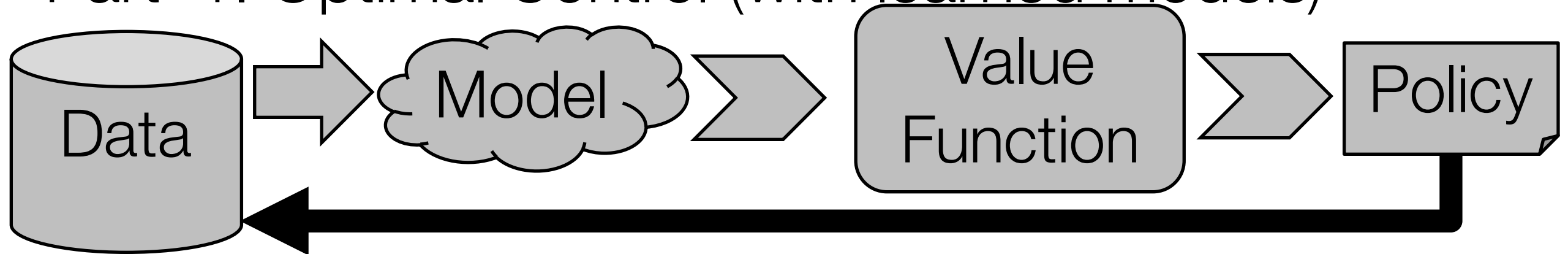
So, why are they not always the method of choice?

- ➔ You need to fill-up you state-action space up with sufficient samples.
- ➔ Another curse of dimensionality with an exponential explosion.
- ➔ Errors in the Value function approximation might have a catastrophic effect on the policy, **can be very hard to control**
- ➔ However, it scales better as we only need samples at relevant locations.

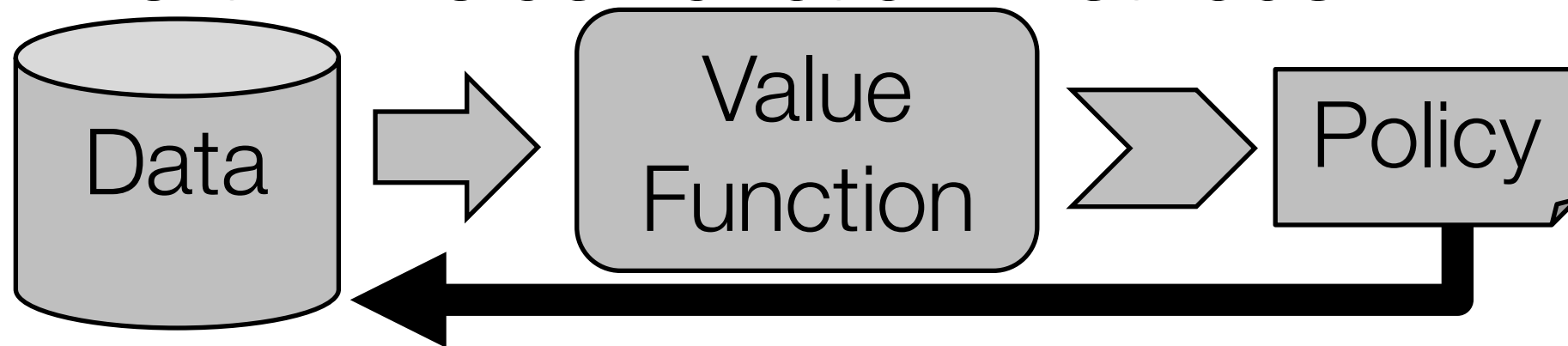
Reinforcement Learning



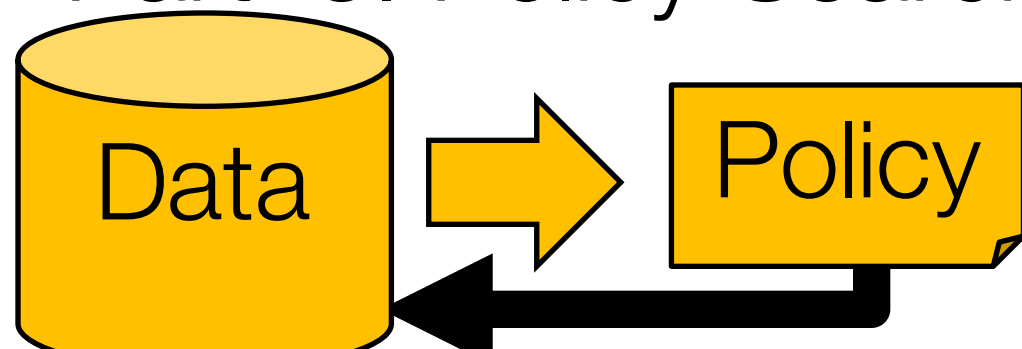
Part 1. Optimal Control (with learned models)



Part 2. Value Function Methods



Part 3. Policy Search



Greedy vs Incremental



Greedy Updates:

$$\theta_{\pi'} = \operatorname{argmax}_{\tilde{\theta}} E_{\pi_{\tilde{\theta}}} \{Q^{\pi}(x, u)\}$$



**potentially
unstable learning
process with large
policy jumps**

Policy Gradient Updates:

$$\theta_{\pi'} = \theta_{\pi} + \alpha \left. \frac{dJ(\theta)}{d\theta} \right|_{\theta=\theta_{\pi}}$$

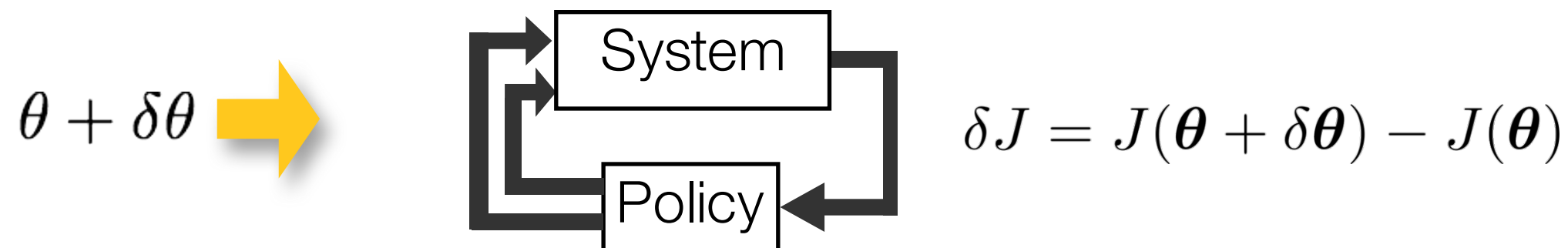


**stable learning
process with
smooth policy
improvement**

Black-Box Approaches, e.g., Finite Differences



1. Perturb the parameters of your policy:



2. Approximate J by first order Taylor approximation

$$J(\theta + \delta\theta) = J(\theta) + \frac{\partial J(\theta)}{\partial \theta} \delta\theta$$

3. Solve for $\frac{\partial J(\theta)}{\partial \theta}$ in a least squares sense (linear regression):

$$\nabla_{\theta}^{\text{FD}} J = \frac{\partial J(\theta)}{\partial \theta} = (\Delta \Theta^T \Delta \Theta)^{-1} \Delta \Theta^T \Delta J$$

Likelihood-Ratio Policy Gradient methods

Some more basic notation

Trajectory distribution: $p(\boldsymbol{\tau}; \boldsymbol{\theta}) = p(\boldsymbol{s}_1) \prod_{t=1}^{T-1} \pi(\boldsymbol{a}_t | \boldsymbol{s}_t; \boldsymbol{\theta}) p(\boldsymbol{s}_{t+1} | \boldsymbol{s}_t, \boldsymbol{a}_t)$

Return for a single trajectory: $R(\boldsymbol{\tau}) = \sum_{t=1}^{T-1} r_t + r_T$

Expected long term reward $J(\boldsymbol{\theta})$ can be written as **expectation over the trajectory distribution**

$$J(\boldsymbol{\theta}) = \mathbb{E}_{p(\boldsymbol{\tau}; \boldsymbol{\theta})}[R(\boldsymbol{\tau})] = \int p(\boldsymbol{\tau}; \boldsymbol{\theta}) R(\boldsymbol{\tau}) d\boldsymbol{\tau}$$



Likelihood Ratio Gradient

The step-based policy gradient can be computed efficiently **by the likelihood-ratio trick**

$$\nabla \log f(x) = \frac{1}{f(x)} \nabla f(x) \quad \Rightarrow \quad \nabla f(x) = f(x) \nabla \log f(x)$$

Applied to the policy gradient

$$\nabla_{\theta} J = \nabla_{\theta} \int p(\tau; \theta) R(\tau) d\tau = \int \nabla_{\theta} p(\tau; \theta) R(\tau) d\tau$$

$$= \int p(\tau; \theta) \nabla_{\theta} \log p(\tau; \theta) R(\tau) d\tau$$

$$\approx \sum_{i=1}^N \nabla_{\theta} \log p(\tau^{[i]}; \theta) R(\tau^{[i]})$$

Needs
only
samples!



Likelihood Ratio Gradient

How do we compute $\nabla_{\boldsymbol{\theta}} \log p(\boldsymbol{\tau}^{[i]}; \boldsymbol{\theta})$?

$$p(\boldsymbol{\tau}; \boldsymbol{\theta}) = p(\mathbf{s}_1) \prod_{t=1}^{T-1} \pi(\mathbf{a}_t | \mathbf{s}_t; \boldsymbol{\theta}) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

The good old log again...

$$\log p(\boldsymbol{\tau}; \boldsymbol{\theta}) = \sum_{t=1}^{T-1} \log \pi(\mathbf{a}_t | \mathbf{s}_t; \boldsymbol{\theta}) + \text{const}$$

Derivative is now easy...

$$\nabla_{\boldsymbol{\theta}} \log p(\boldsymbol{\tau}; \boldsymbol{\theta}) = \sum_{t=1}^{T-1} \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}_t | \mathbf{s}_t; \boldsymbol{\theta})$$

Lets plug it in...



Result:

$$\begin{aligned}\nabla_{\theta} J &= \sum_{i=1}^N \sum_{t=1}^{T-1} \nabla_{\theta} \log \pi(\mathbf{a}_t^{[i]} | \mathbf{s}_t^{[i]}; \theta) R(\tau^{[i]}) \\ &= \sum_{i=1}^N \sum_{t=1}^{T-1} \nabla_{\theta} \log \pi(\mathbf{a}_t^{[i]} | \mathbf{s}_t^{[i]}; \theta) \left(\sum_{t=1}^{T-1} r_t^{[i]} + r_T^{[i]} \right)\end{aligned}$$

This algorithm is called the **REINFORCE Policy Gradient**

Does this method work well?

No!



Kullback Leibler divergences

The Natural gradient is defined as the update direction which is closest to the standard gradient, **but has limited distance to the old distribution**

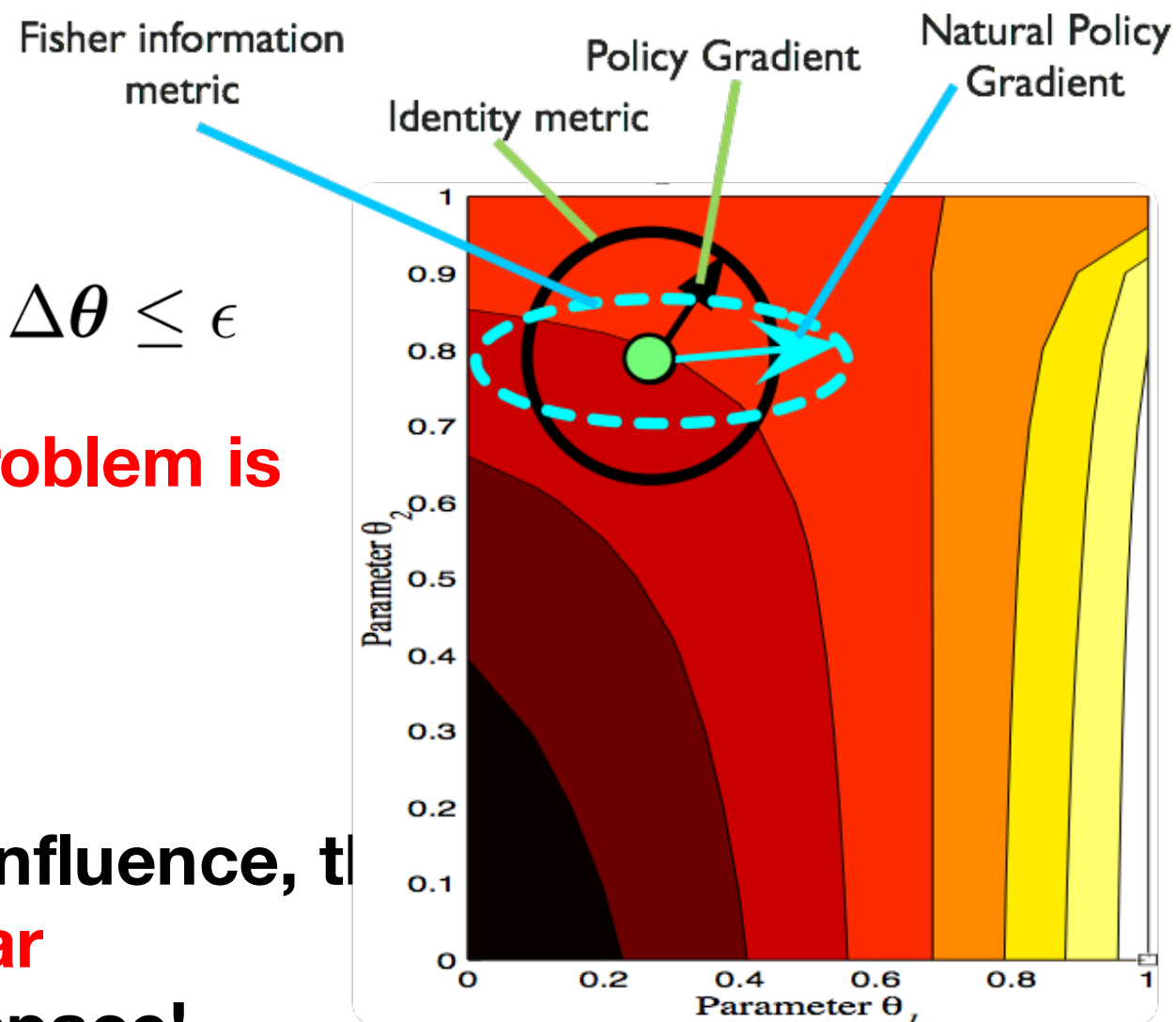
$$\nabla_{\theta}^{\text{NG}} J = \operatorname{argmax}_{\Delta\theta} \Delta\theta^T \nabla_{\theta} J$$

$$\text{s.t.: } \text{KL}(p_{\theta+\Delta\theta} || p_{\theta}) \approx \Delta\theta^T G(\theta) \Delta\theta \leq \epsilon$$

The solution to this optimization problem is given as:

$$\nabla_{\theta}^{\text{NG}} J \propto G(\theta)^{-1} \nabla_{\theta} J$$

As every parameter has the same influence, the natural gradient is **invariant to linear transformations** of the parameter space!



Are they useful?

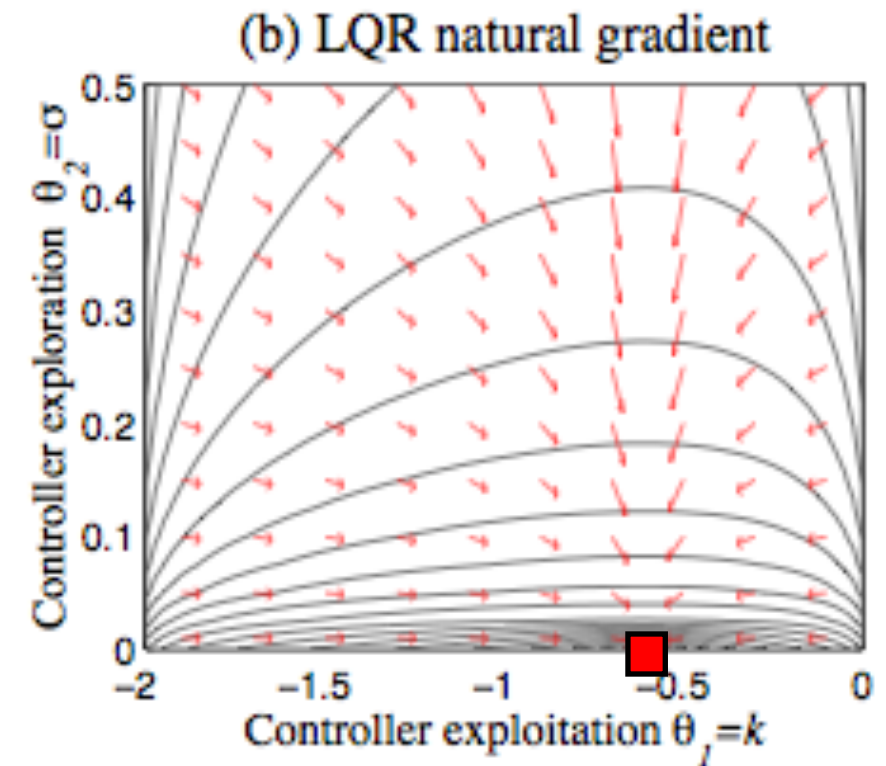
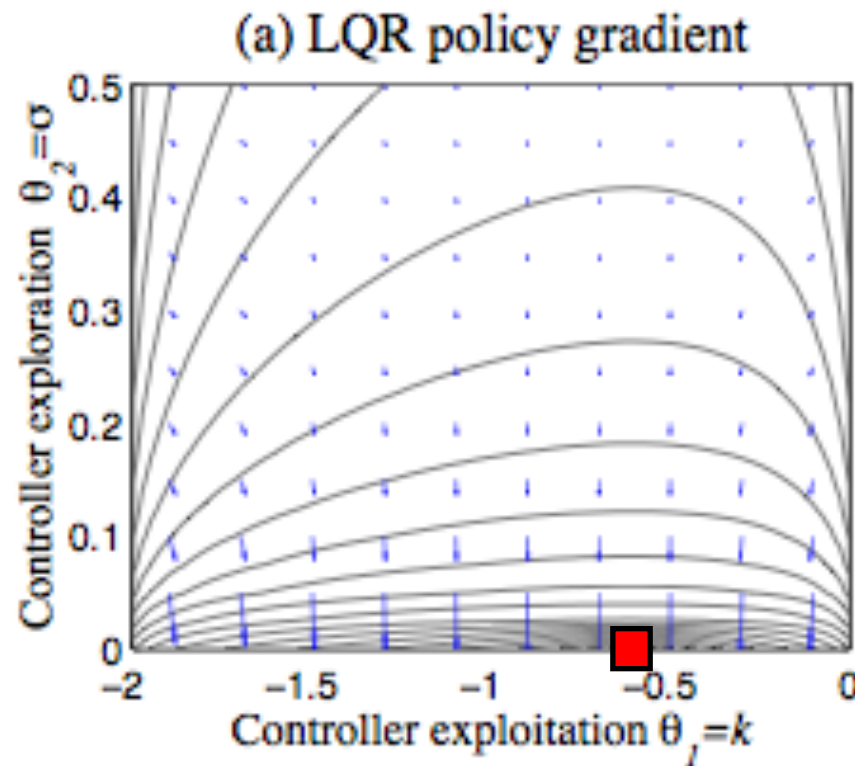


Linear Quadratic Regulation

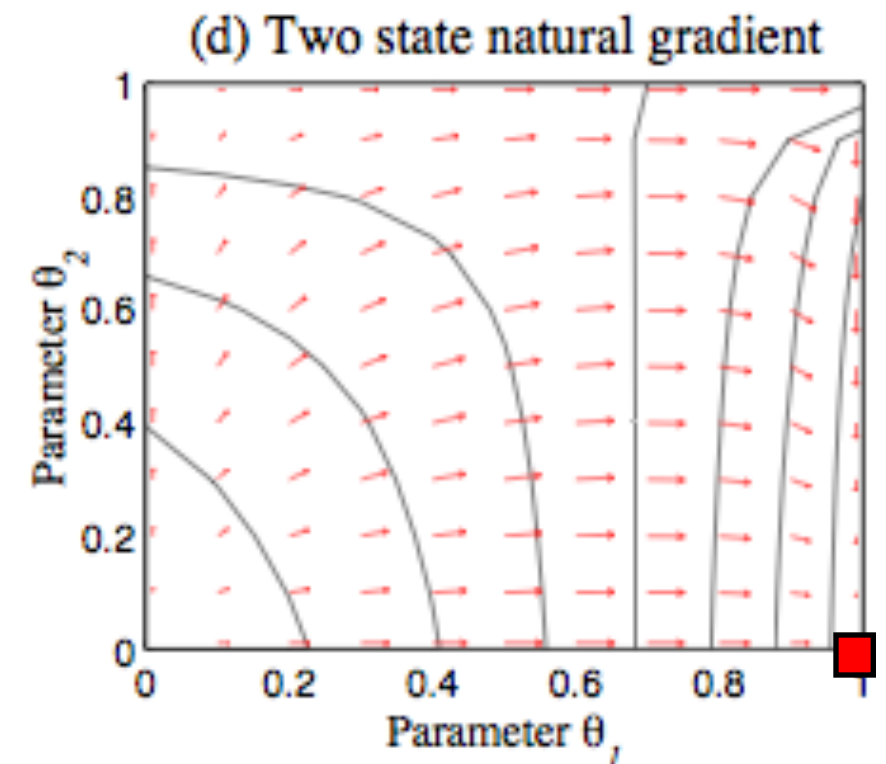
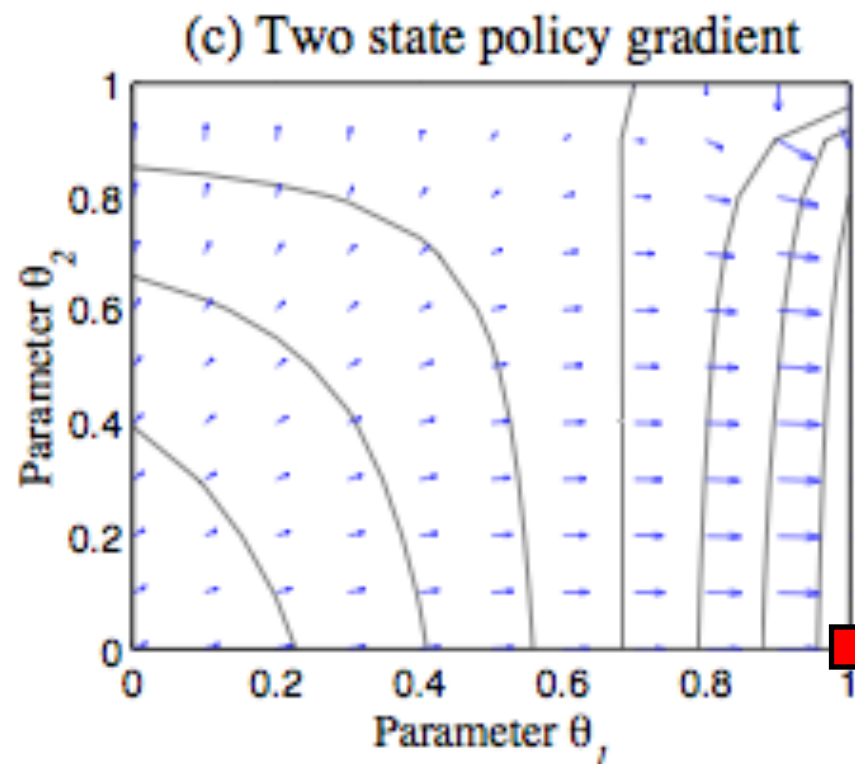
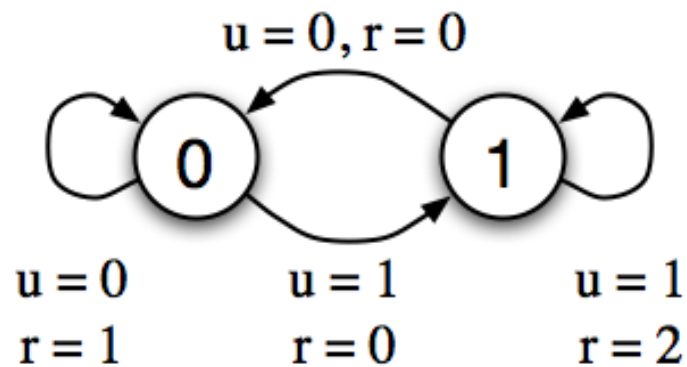
$$x_{t+1} = Ax_t + Bu_t$$

$$u_t \sim \pi(u|x_t) = \mathcal{N}(u|kx_t, \sigma)$$

$$r_t = -x_t^T Q x_t - u_t^T R u_t$$



Two-State Problem



(Peters et al. 2003, 2005)

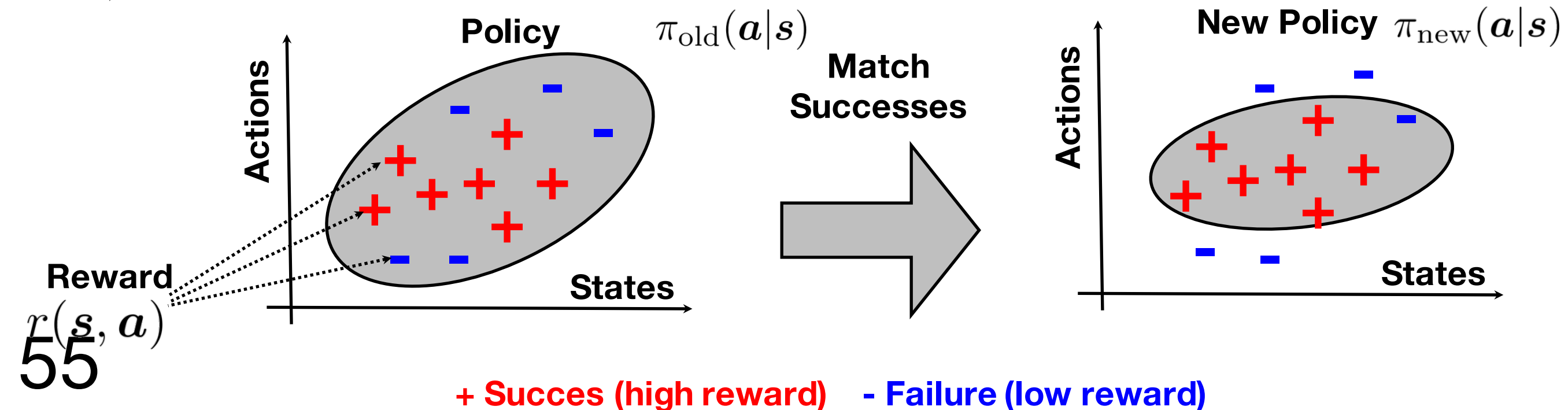


Success Matching Principle

“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 $r(\mathbf{s}, \mathbf{a})$



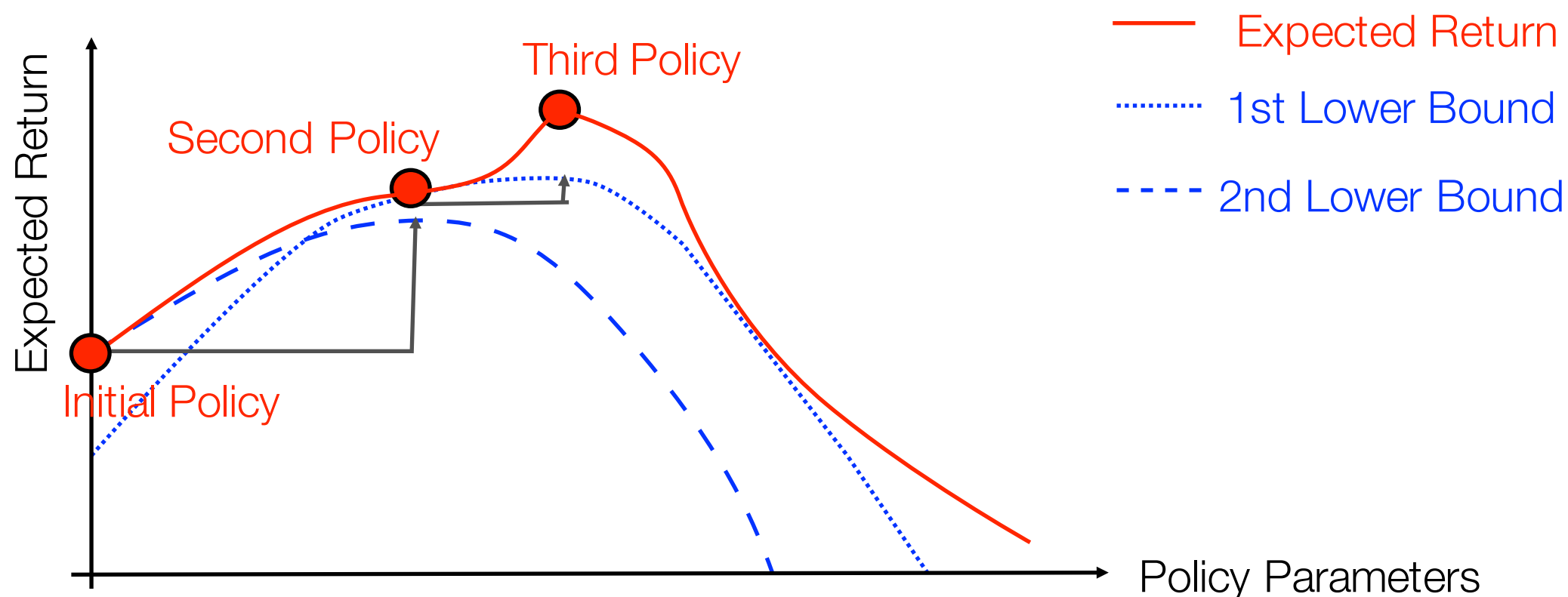
Basic Intuition



- Lower Bound on Expected Return
 - reward is an improper probability distribution
 - log-likelihood \rightarrow log(expected return)

(Dayan & Hinton, Neural Computation 1997; Peters & Schaal, ICML 2007; Kober & Peters, NIPS 2008)

$$\log J(\theta') \geq \int_{\mathbb{T}} p_{\theta}(\tau) R(\tau) \log \frac{p_{\theta'}(\tau)}{p_{\theta}(\tau)} d\tau + \text{const} = L_{\theta}(\theta')$$



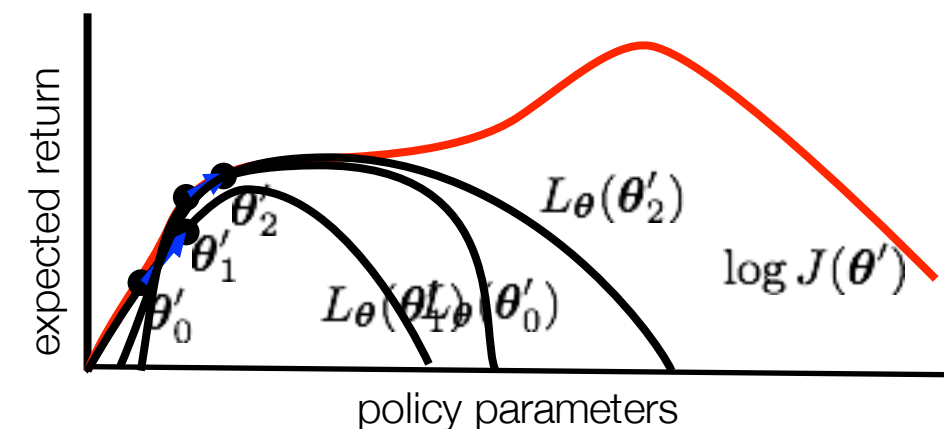
Resulting Algorithms



Policy Gradients: maximize lower bound by following the gradient

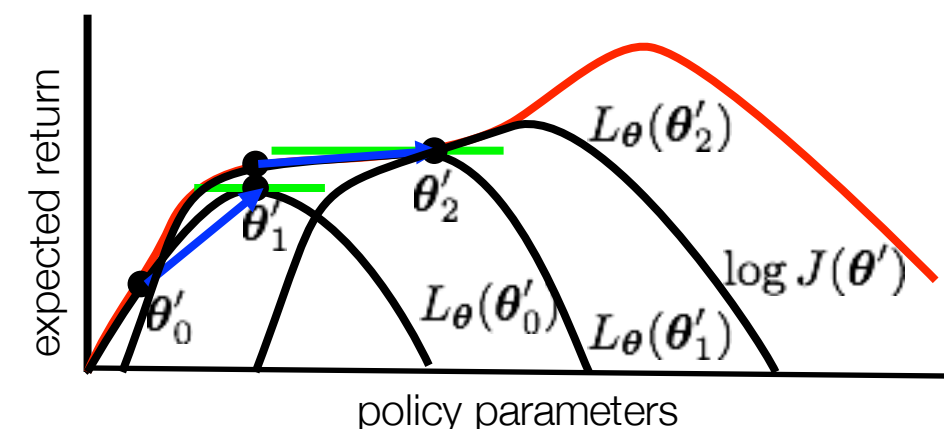
$$\lim_{\theta' \rightarrow \theta} \partial_{\theta'} L_{\theta}(\theta') = \partial_{\theta} J(\theta)$$

$$\theta' \approx \theta + \alpha \partial_{\theta} J(\theta)$$



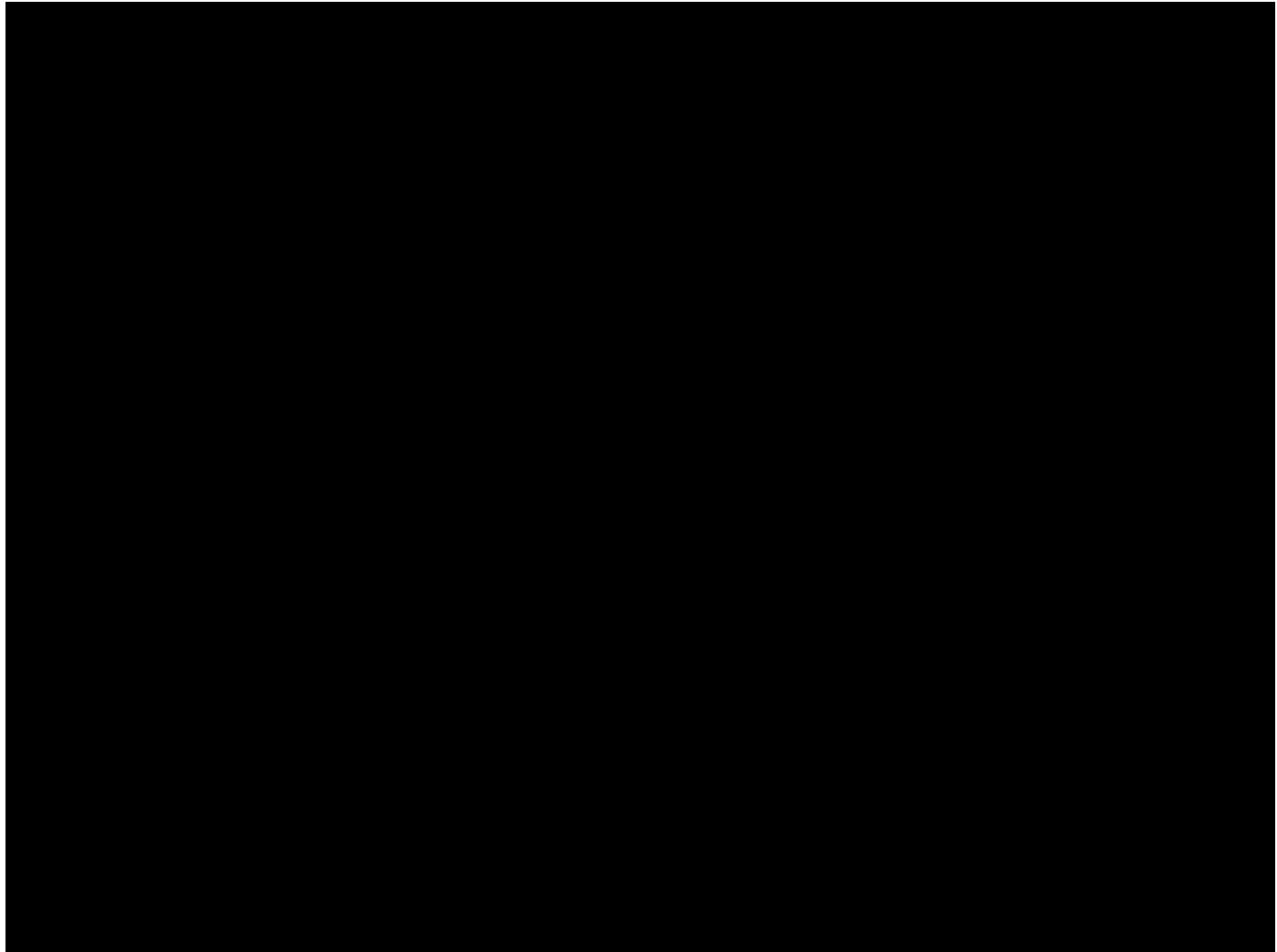
EM-like Methods: maximize lower bound by expectation-maximization

$$\theta' = \operatorname{argmax} L_{\theta}(\theta')$$





Ball in the Cup



What are the key problems?

Key Problems

1.no notion of data in the generic problem formulation

Let's introduce the *observed data distribution* $q(\mathbf{x}, \mathbf{u})$.

2.optimization bias problematic with data

Let's bound our information loss $D(\mu^{\pi_\theta}(\mathbf{x})\pi_\theta(\mathbf{u}|\mathbf{x})||q(\mathbf{x}, \mathbf{u})) \leq \epsilon$

3.role of features is unclear in most methods

Let's introduce stationary features determined by

$$\sum_{\mathbf{x}'} \mu^\pi(\mathbf{x}') \phi_{\mathbf{x}'} = \sum_{\mathbf{x}, \mathbf{u}, \mathbf{x}'} \mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x}) p(\mathbf{x}'|\mathbf{x}, \mathbf{u}) \phi_{\mathbf{x}'}$$

Relative Entropy Policy Search Problem

Peters et al., AAAI 2010

$$\max_{\mu^\pi, \pi} J(\pi) = \sum_{\mathbf{x}, \mathbf{u}} \mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x}) r(\mathbf{x}, \mathbf{u})$$

s.t.

$$\sum_{\mathbf{x}'} \mu^{\pi'}(\mathbf{x}') \phi(\mathbf{x}') = \sum_{\mathbf{x}, \mathbf{u}, \mathbf{x}'} \mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x}) p(\mathbf{x}'|\mathbf{x}, \mathbf{u}) \phi(\mathbf{x}')$$

$$1 = \sum_{\mathbf{x}, \mathbf{u}} \mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x})$$

$$\sum_{\mathbf{x}, \mathbf{u}} \mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x}) \log \frac{\mu^\pi(\mathbf{x}) \pi(\mathbf{u}|\mathbf{x})}{q(\mathbf{x}, \mathbf{u})} \leq \epsilon$$

Solution

Peters et al., AAAI 2010

New Gibbs Policy

Value Function

$$\pi(\mathbf{u}|\mathbf{x}) = \frac{q(\mathbf{x}, \mathbf{u}) \exp\left(\frac{1}{\eta} \delta_{\theta}(\mathbf{x}, \mathbf{u})\right)}{\sum_{\mathbf{a}} q(\mathbf{x}, \mathbf{a}) \exp\left(\frac{1}{\eta} \delta_{\theta}(\mathbf{x}, \mathbf{a})\right)}$$

Advantage Function given as $V_{\theta}(\mathbf{x}) = \phi_{\mathbf{x}}^T \theta$

Parameters

$$\delta_{\theta}(\mathbf{x}, \mathbf{u}) = r(\mathbf{x}, \mathbf{u}) + \sum_{\mathbf{x}'} p(\mathbf{x}'|\mathbf{x}, \mathbf{u}) V_{\theta}(\mathbf{x}') - V_{\theta}(\mathbf{x})$$

$$\max_{\theta, \eta} g(\theta, \eta) = \eta \log \left(\sum_{\mathbf{x}, \mathbf{u}} q(\mathbf{x}, \mathbf{u}) \exp \left(\epsilon + \frac{1}{\eta} \delta_{\theta}(\mathbf{x}, \mathbf{u}) \right) \right)$$

All direct results of the previous problem!



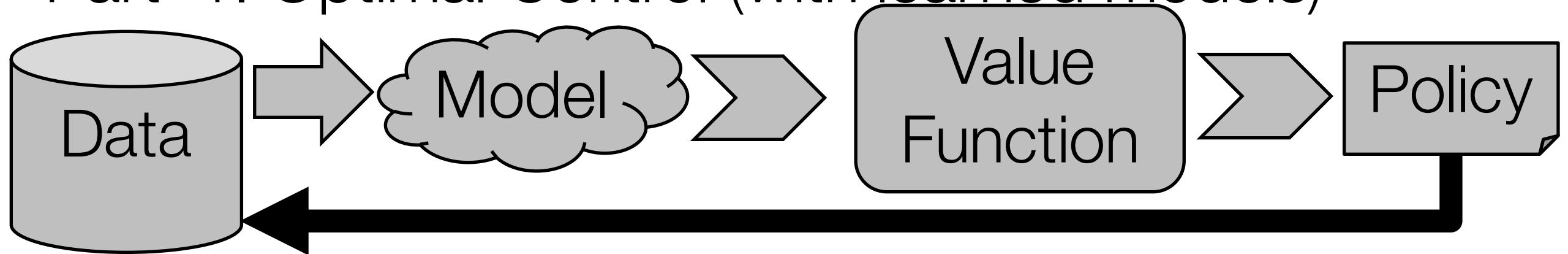
Wrap-Up

- 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.
- Learning the exploration rate is still an open problem
- Newer methods focus on probabilistic policy search approaches.
- Relative Entropy Policy Search (and its simplifications like TRPO) would be today's choice!

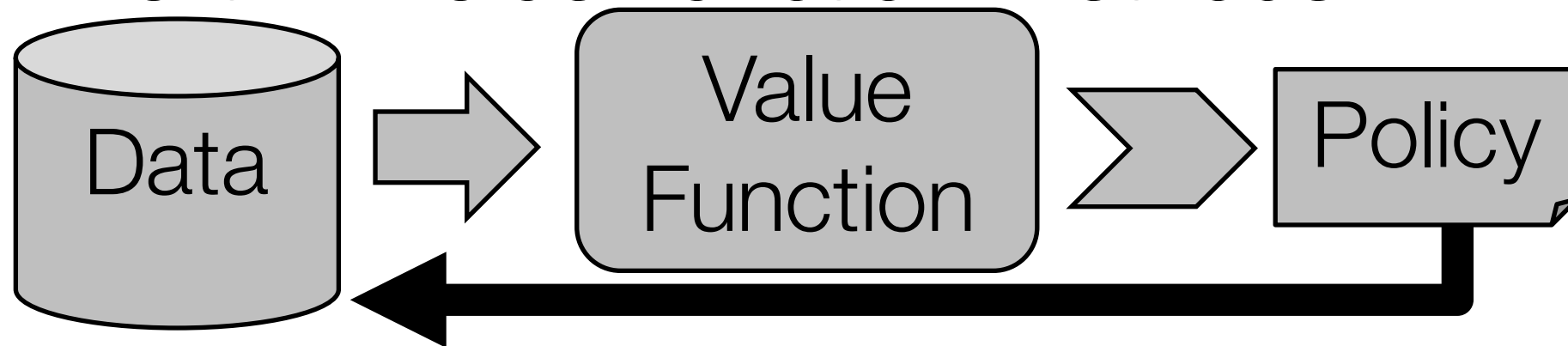
Reinforcement Learning



Part 1. Optimal Control (with learned models)



Part 2. Value Function Methods



Part 3. Policy Search

