

Trade-Offs in Robot Skill Learning

Jan Peters Technische Universität Darmstadt

Max Planck Institute for Intelligent Systems



TECHNISCHE UNIVERSITÄT DARMSTADT





Motivation

Can we create such robots?

Source: Movie iRobot



Motivation



Adapt to humans



Uncertainty in tasks and environment

Programming complexity beyond human imagination

How can we fulfill Hollywood's vision of future robots?

- Smart Humans? Hand-coding of behaviors has allowed us to go very far!
- Maybe we should allow the robot to learn new tricks, adapt to situations, refine skills?
- Why have "Off-the-shelf" machine learning approaches not delivered the need intelligence?

We trade-off between autonomy, human insight and behavior quality to get close!

Trade-Offs in Robot Skill Learning

Behavior Quality Efficiency of Generation

Trade-Offs in Robot Learning

Autonomy & Generality



Outline

State

Task Parameters

- I. Introduction: Robot Skill Learning in Four Slides
- 2. Trade-Off I: Generality vs Physical Knowledge
- 3. Trade-Off II: Information Loss vs Self-Improvement
- 4. Trade-Off III: Large Repertoire vs Parsimonious Behavior Representation





Example



Internal and external state x_t , action u_t .

Modeling Assumptions



Policy: Generates action \mathbf{u}_t in state \mathbf{x}_t . Should we use a deterministic policy $\mathbf{u}_t = \pi(\mathbf{x}_t)$? **NO!** Stochasticity is important: **Robot learning**

- needed for exploration
- breaks "curse of dimensionality"
- optimal solution can be stochastic

Hence, we use a stochastic policy: $\mathbf{u}_t \sim \pi(\mathbf{u}_t | \mathbf{x}_t)$

Teacher: Evaluates the performance and rates it with r_t .

Environment: An action \mathbf{u}_t causes the system to change state from \mathbf{x}_t to \mathbf{x}_{t+1} . Model in the real world: $\mathbf{x}_{t+1} \sim p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t)$

implies "policy

optimization"!

Let the loop roll out!



Trajectories

$$\boldsymbol{ au} = [\mathbf{x}_0, \mathbf{u}_0, \mathbf{x}_1, \mathbf{u}_1 \dots, \mathbf{x}_{T-1}, \mathbf{u}_{T-1}, \mathbf{x}_T]$$

Path distributions

$$p(\tau) = p(\mathbf{x}_0) \prod_{t=0}^{T-1} p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) \pi(\mathbf{u}_t | \mathbf{x}_t)$$

Path rewards:

$$r(\boldsymbol{\tau}) = \sum_{t=0}^{T} \alpha_t r(\mathbf{x}_t, \mathbf{u}_t)$$

What is robot learning?

In our model: Optimize the expected scores $J(\theta) = E_{\tau}\{r(\tau)\} = \int_{\mathbb{T}} p_{\theta}(\tau)r(\tau)d\tau$ of the teacher.

Peters & Schaal (2003). Reinforcement Learning for Humanoid Robotics, HUMANOIDS



Outline

State

Task Parameters

- I. Introduction: Robot Skill Learning in Four Slides
- 2. Trade-Off I: Generality vs Physical Knowledge
- 3. Trade-Off II: Information Loss vs Self-Improvement
- 4. Trade-Off III: Large Repertoire vs Parsimonious Behavior Representation









- Physical knowledge is extremely powerful!
- Physics yields general solutions!
- Physics can explain itself!
- Handcrafting with Physics Knowledge can be easy



- Physical knowledge is extremely powerful, but ...
- ...classical robotics as well as classical AI have failed at encoding the whole world with unlimited precision.
- ...all physical models are wrong (but some are useful)
- learning allows for higher accuracy, robustness and autonomy!

Trade-Off I: Physical Knowledge vs Generality

Optimal Trade-Off

Learning beats physical knowledge if we extract crucial features...

Physical Model

Use all physics...



Offline Trained Model

Online Trained

Learn everything...



Lesson: We need a physics-based blue print!

Nguyen-Tuong & Peters, ARJ 2010 14



A Physics-inspired Blue Print Optimal Trade-Off Use all physics... Learn everything... Task/Hyperparameter $\dot{z} = \alpha_z (\beta_z (g - y) - z)$ $\dot{y} = \alpha_y (f(x, v) + z)$ Trajectory Plan **Dynamics** where Linear in learnable $\begin{cases} \dot{v} = \alpha_v (\beta_v (g - x) - v) \\ \end{cases} Policy Parameters \end{cases}$ Canonical **Dynamics** $\dot{x} = \alpha_v v$ $f(x,v) = \frac{\sum_{i=1}^{n} w_i b_i v_i}{\sum_{i=1}^{k} w_i}$ Local Linear Model Approx. $w_i = \exp\left(-\frac{1}{2}d_i(\overline{x}-c_i)^2\right)$ and $\overline{x} = \frac{x-x_0}{g-x_0}$

Imitation Learning

Given a physics-based policy, can we reproduce a path distribution?

• match given path distribution p(T) with a new one $p_{\theta}(T)$, i.e.,

 $D(p_{\theta}(\boldsymbol{\tau})||p(\boldsymbol{\tau})) \to \min$

- only adapt the policy parameters θ
- model-free, purely samplebased
- results in one-shot and expectation maximization algorithms





Acquisition by Imitation

Trade-Offs in Robot Skill Learning

Behavior Quality

Online Learning for higher Accuracy Physics-knowledge in the representation Efficiency of Generation

Trade-Offs in Robot Learning

Sufficiently general representations

Autonomy & Generality



Outline

State

Task Parameters

- I. Introduction: Robot Skill Learning in Four Slides
- 2. Trade-Off I: Generality vs Physical Knowledge
- 3. Trade-Off II: Information Loss vs Self-Improvement
- 4. Trade-Off III: Large Repertoire vs Parsimonious Behavior Representation



Trade-Off II: Information Loss vs Self-Improvement



Reinforcement Learning

Given a path distribution, can we find the optimal policy?

- Goal: maximize the return of the paths r(T) generated by path distribution $p_{\theta}(T)$
- Optimization function is the expected reward

$$J(\boldsymbol{\theta}) = \int_{\mathbb{T}} p_{\boldsymbol{\theta}}(\boldsymbol{\tau}) r(\boldsymbol{\tau}) d\boldsymbol{\tau}$$

- This optimal control problem has no data in it!
- This part usually results into a greedy, softmax updates or a `vanilla' policy gradient algorithm...

Success Matching

"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).

Can we create better policies by matching the rewardweighted previous policy ?



Many related frameworks, e.g., (Dayan&Hinton 1992;Andrews,'03;Attias,'04;Bagnell,'03;Toussaint,'06;...). 23

Illustrative Example Foothold Selection



Match successful footholds!

What about the new cost function?

Does the new goal function help?

• We have a lower bound

$$\log J(\boldsymbol{\theta}) = \log \int \frac{p(\boldsymbol{\tau})}{p(\boldsymbol{\tau})} p_{\boldsymbol{\theta}}(\boldsymbol{\tau}) r(\boldsymbol{\tau}) d\boldsymbol{\tau},$$

$$\geq \int p(\boldsymbol{\tau}) r(\boldsymbol{\tau}) \log \frac{p_{\boldsymbol{\theta}}(\boldsymbol{\tau})}{p(\boldsymbol{\tau})} d\boldsymbol{\tau} \propto -D(p_{\boldsymbol{\theta}}(\boldsymbol{\tau})||p(\boldsymbol{\tau})r(\boldsymbol{\tau})),$$

- Having taken a low reward trajectory many times resembles having taken a high reward trajectory once!
- A safe way of trading experience against reward!



Reinforcement Learning by Reward-Weighted Imitation

Matching successful actions corresponds to minimizing the Kullback-Leibler 'distance'

 $D(p_{\theta}(\boldsymbol{\tau})||r(\boldsymbol{\tau})p(\boldsymbol{\tau})) \to \min$

For a Gaussian policy $\pi(\mathbf{u}|\mathbf{x}) = \mathcal{N}(\mathbf{u}|\boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\theta}, \sigma^2 \mathbf{I})$, we get the update rule



Reduces Reinforcement Learning onto Reward Weighted Regression!

Peters & Schaal (2007). Policy Learning for Motor Skills, International Conference on Machine Learning (ICML) Kober & Peters (2009). Policy Search for Motor Primitives in Robotics, Advances in Neural Information Processing Systems (NIPS)



Self-Improvement by Reinforcement Learning

27

More focussed trade-off?

EM-like Methods: Safe but unclear Trade-Off Experience High Reward REPS: Adjustable Trade-Off

Relative Entropy Policy Search (REPS)

$$\max_{\pi,\mu^{\pi}} J(\pi) = \sum_{s,a} \mu^{\pi}(s) \pi(a|s) \mathcal{R}_{sa} \text{ Maximize reward}$$

$$1 = \sum_{s,a} \mu^{\pi}(s) \pi(a|s)$$
 Probability distribution

$$\mu^{\pi}(s') = \sum_{s,a} \mathcal{P}^{a}_{ss'} \mu^{\pi}(s) \pi(a|s)$$
 Follow system dynamics

$$\epsilon \geq \sum_{s,a} \mu^{\pi}(s) \pi(a|s) \log \frac{\mu^{\pi}(s) \pi(a|s)}{q(s,a)}$$
Close to training data (no wild exploration)

Peters, Muelling, Altun (2010). Relative Entropy Policy Search, AAAI 28

Trade-Offs in Robot Skill Learning

Behavior Quality

Highly rewarded behaviors!

Safe Learning Approaches by Retaining Experience!

Initial Demonstrations

Trade-Offs in Robot Learning Efficiency of Generation

Slow Loss of (Bad) Experience

Autonomy & Generality



Outline

Task Parameters

- I. Introduction: Robot Skill Learning in Four Slides
- 2. Trade-Off I: Generality vs Physical Knowledge
- 3. Trade-Off II: Information Loss vs Self-Improvement
- 4. Trade-Off III: Large Repertoire vs Parsimonious Behavior Representation
- 5. Conclusion Current State Compand



- Versatility: Complex behaviors rely on a combination of as many strategies as possible.
- Parsimony: Multiple, distinct solutions are required for efficient learning.

Versatility with Many Primitives

Relative Entropy Policy Search (REPS)

$$\begin{split} \max_{\pi,\mu^{\pi}} J(\pi) &= \sum_{s,a} \mu^{\pi}(s) \pi(a|s) \mathcal{R}_{sa} \quad \text{Maximize reward} \\ 1 &= \sum_{s,a} \mu^{\pi}(s) \pi(a|s) \quad \text{Probability distribution} \\ \mu^{\pi}(s') &= \sum_{s,a} \mathcal{P}_{ss'}^{a} \mu^{\pi}(s) \pi(a|s) \quad \text{Follow system dynamics} \\ \epsilon &\geq \sum_{s,a} \mu^{\pi}(s) \pi(a|s) \log \frac{\mu^{\pi}(s) \pi(a|s)}{q(s,a)} \quad \text{Close to training data (no wild exploration)} \end{split}$$



Hierarchies with Primitives



Daniel, Neumann & Peters (2012). Hierarchical Relative Entropy Policy Search, AI-Stats 34

Versatility with Many Primitives

"Naive" Extension of REPS

$$\max_{\pi,\mu^{\pi}} J(\pi) = \sum_{s,a,o} \mu^{\pi}(s)\pi(a,o|s)\mathcal{R}_{sa} \text{ Maximize reward}$$

$$1 = \sum_{s,a,o} \mu^{\pi}(s)\pi(a,o|s) \text{ Probability distribution}$$

$$\mu^{\pi}(s') = \sum_{s,a,o} \mathcal{P}_{ss'}^{a}\mu^{\pi}(s)\pi(a|o,s)\pi(o|s) \text{ Follow system dynamics}$$

$$\epsilon \ge \sum_{s,a,o} \mu^{\pi}(s)\pi(a,o|s)\log\frac{\mu^{\pi}(s)\pi(a,o|s)}{q(s,a,o)} \text{ Close to training data (no wild exploration)}$$

Daniel, Neumann & Peters (2012). Hierarchical Relative Entropy Policy Search, Al-Stats 35



If all primitives are equally responsible, we can represent versatile behavior but it will never be parsimonious.

Daniel, Neumann & Peters (2012). Hierarchical Relative Entropy Policy Search, AI-Stats

Versatility with Many Primitives

Hierachical Relative Entropy Policy Search (HiREPS)

 $\max_{\pi,\mu^{\pi}} J(\pi) = \sum_{s,a,o} \mu^{\pi}(s) \pi(a,o|s) \mathcal{R}_{sa} \text{ Maximize reward}$

 $1 = \sum_{s,a,o} \mu^{\pi}(s) \pi(a,o|s)$ Probability distribution

 $\mu^{\pi}(s') = \sum_{s,a,o} \mathcal{P}^{a}_{ss'} \mu^{\pi}(s) \pi(a|o,s) \pi(o|s)$ Follow system dynamics

High entropy indicates overlap!

$$\epsilon \geq \sum_{s,a,o} \mu^{\pi}(s) \pi(a,o|s) \log \frac{\mu^{\pi}(s)\pi(a,o|s)}{q(s,a,o)} \text{ Close to training data (no wild exploration)}$$
$$\kappa \geq \mathbb{E}_{s,a} \left[\sum_{s,a} -p(o|s,a) \log p(o|s,a) \right] \text{ Force the primitives to limited responsibility}$$

ο



We can reduce to the number of needed primitives!

Daniel, Neumann & Peters (2012). Hierarchical Relative Entropy Policy Search, Al-Stats 38



Pole

Target Zone

2 Options

Daniel, Neumann & Peters (2012). Hierarchical Relative Entropy Policy Search, Al-Stats

String



Good performance

Fast reduction in the number of primitives

Daniel, Neumann & Peters (2012). Hierarchical Relative Entropy Policy Search, Al-Stats



Demonstrations

Demonstrations with Kinesthetic Teach-In

Select & Generalize

From Imitation Learning we obtain 25 Movement Primitives

Self-Improvement

Training a Hitting Region with an Initial Success Rate of 0%

Current Gameplay

Final Challenge: Match against a Human

Trade-Offs in Robot Skill Learning





Outline

Task Parameters

- I. Introduction: Robot Skill Learning in Four Slides
- 2. Trade-Off I: Generality vs Physical Knowledge
- 3. Trade-Off II: Information Loss vs Self-Improvement
- 4. Trade-Off III: Large Repertoire vs Parsimonious Behavior Representation

State Conclusion Execute Action Action

Conclusion

- Motor skill learning is a promising way to avoid programming all possible scenarios and continuously adapt to the environment.
- To make it work, we always have to trade-off at least:
 - Physical Knowledge vs Learning
 - Experienced vs High Rewards
 - Versatility vs Parsimony
- There are many more trade-offs in practice.

Thanks for your Attention!





Boularias



Amor



Gerhard Neumann



Roberto Calandra





Herke van Hoof

Alexandros Paraschos





Marc Deisenroth

> Christian Daniel

