Generalization and Safety in Reinforcement Learning and Control

Aviv Tamar

Berkeley

Artificial Intelligence Research Laboratory

Today: robots in factories



- Future: robots 'in the wild'
 - Inside homes
 - Alongside humans
 - Rapid manufacturing







- Future: robots 'in the wild'
 - Inside homes
 - Alongside humans
 - Rapid manufacturing
- Challenges:
 - Unstructured environments
 - Nonlinear, non-smooth dynamics
 - True model is never known
 - Safety is essential





mar, UC Berkeley

- Future: robots 'in the wild'
 - Inside homes
 - Alongside humans
 - Rapid manufacturing
- Challenges:
 - Unstructured environments
 - Nonlinear, non-smooth dynamics
 - True model is never known
 - Safety is essential





mar, UC Berkeley

Reinforcement Learning

Reinforcement learning (RL) model:



Reinforcement Learning

Reinforcement learning (RL) model:



 Machine learning to learn (state) → (action) that leads to high reward

Reinforcement Learning

- Naturally addresses:
 Unstructured environments
 Nonlinear, non-smooth dynamics
 True model is never known
- Practical challenges:
 - Safety
 - Generalization (to changes in the task)
 - High-dimensional states/actions







This Talk

- 1. Safety risk sensitive RL
- 2. Generalization value iteration networks
- 3. Planning-based neural networks for robotics

Part 1: Risk-Aware RL

Background

Reinforcement Learning

MDP Formulation

- States s, actions a
- Transitions P(s'|s, a)
- Policy $\pi(a|s)$



- Reward r(s)
- Discount $\gamma \in [0, 1)$

• Goal:
$$\max_{\pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) | \pi \right]$$

Background

Planning (when model is known)

• Value function:

$$V^{\pi}(s) = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^{t} r(s_{t}) | \pi, s_{0} = s\right]$$

• Value iteration algorithm:

$$V_{n+1}(s) = \max_{a} Q_n(s, a) \quad \forall s,$$
$$Q_n(s, a) = R(s) + \gamma \sum_{s'} P(s'|s, a) V_n(s').$$

Background

Planning (when model is known)

• Value function:

$$V^{\pi}(s) = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^{t} r(s_{t}) | \pi, s_{0} = s\right]$$

• Value iteration algorithm:

$$V_{n+1}(s) = \max_{a} Q_n(s, a) \quad \forall s,$$
$$Q_n(s, a) = R(s) + \gamma \sum_{s'} P(s'|s, a) V_n(s').$$

- Converges to $V^* = \max_{\pi} V^{\pi}$ (γ -geometric rate)
- Optimal policy $\pi^*(a|s) = \arg \max_a Q^*(s, a)$

RL (when model unknown/too large to plan)

- 1. Approx. value function (e.g., neural-net): $V^*(s) \approx V(s;\theta)$
 - Sampling to estimate θ from interaction
- 2. Policy gradient:

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

$$\nabla_{\theta} \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^{t} r(s_{t}) | \pi\right]$$

• Sampling to estimate gradient from interaction

Risk Averse RL

- Why should we care about risk?
 - 1. The world is stochastic
 - Risk-aware goal: $\max_{\pi} \rho \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) | \pi \right]$
 - ρ is risk measure, e.g., $\rho(X) = \mathbb{E}[X] \beta \operatorname{Var}[X]$

Safety against unluckiness



Risk Averse RL

- Why should we care about risk?
 - 2. True world model is unknown
 - Robustness objective:

 $\max_{\pi} \min_{\hat{P} \in \text{ possible MDPs}} \mathbb{E} \left| \sum_{t=0}^{\infty} \gamma^{t} r(s_{t}) \right| \pi, \hat{P} \right|$

Safety against model mismatch



Risk Averse RL

- Why should we care about risk?
 - 1. Safety against unluckiness
 - 2. Safety against model mismatch





Conditional Value at Risk (CVaR)

- Safety against 'unluckiness'
- Definition:
 - X random variable
 - $q_{\alpha}(X)$ α quantile
 - $\operatorname{CVaR}_{\alpha}(X) = \mathbb{E}[X|X \leq q_{\alpha}]$
 - Expected α % worst cases



Risk averse goal:

$$\max_{\pi} \mathbf{CVaR} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) | \pi \right]$$

Conditional Value at Risk (CVaR)

- Popular in finance
- CVaR Risk = Robustness^[1]
 - Multiplicative perturbations

$$\hat{P}(x_{t+1}|x_t) = P(x_{t+1}|x_t) \cdot \delta_t(x_{t+1}|x_t)$$

• Budget

 $\delta_1(x_1|x_0)\delta_2(x_2|x_1)\cdots\delta_H(x_H|x_{H-1}) \leq \eta, \quad \forall x_1,\ldots,x_H.$

The worst cannot happen at every time!

• Set of possible perturbations Δ_{η}

Theorem

$$CVaR_{\frac{1}{\eta}}\left(\sum_{t=0}^{H} r(x_t)\right) = \inf_{\hat{P} \in \Delta_{\eta}} \mathbb{E}\left[\left|\sum_{t=0}^{H} r(x_t)\right| \hat{P}\right]$$

[1] Chow, T., Mannor, Pavone. NIPS 2015



Conditional Value at Risk (CVaR)

Safety against unluckiness = safety to model errors

- CVaR Can be optimized efficiently:
 - 1. Planning^[1]
 - Value iteration + linear prog. + linear interpolation
 - First error bounds + convergence rate
 - γ -geometric rate, poly(S, A) each iteration
 - 2. Large/continuous MDPs, RL^[2,3]
 - Policy gradient
 - Convergence (w.p.1) to locally optimal policy
 - Compatible with deep RL
- [1] Chow, T., Mannor, Pavone. NIPS 2015
- [2] T., Glassner, Mannor. AAAI 2015
- [3] T., Chow, Ghavamzadeh, Mannor. NIPS 2015

CVaR definition

- X random variable
- $q_{\alpha}(X)$ α quantile
- α -CVaR:

 $\Phi_{\alpha}(X) = \mathbb{E}\left[\left. X \right| X \leq q_{\alpha} \right]$

- Expected α % worst cases
- Prominent in finance
- Sensitive to rare, disastrous events





T., Glassner, Mannor. AAAI 2015

Gradient estimation

- Likelihood ratio method (Glynn 1990; a.k.a. policy gradient)
- Estimate gradient $\nabla \mathbb{E}(X)$

 ∇

$$E(X) = \nabla \int_{-\infty}^{\infty} f_X(x) x dx$$

= $\int_{-\infty}^{\infty} \nabla f_X(x) x dx$
= $\int_{-\infty}^{\infty} \frac{\nabla f_X(x)}{f_X(x)} f_X(x) x dx$
= $\mathbb{E}\left(\frac{\nabla f_X(X)}{f_X(X)} X\right)$
 $\approx \frac{1}{N} \sum_{1 \le i \le N} \frac{\nabla f_X(x_i)}{f_X(x_i)} x_i$

Gradient estimation - CVaR

• Estimate gradient $\nabla \Phi_{\alpha}(X)$

Maybe:

$$\nabla \Phi_{\alpha}(X) \approx \frac{1}{\alpha N} \sum_{\alpha N \text{ worst}} \frac{\nabla f_X(x_i)}{f_X(x_i)} x_i$$

Gradient estimation - CVaR

- Estimate gradient $\nabla \Phi_{\alpha}(X)$
- Maybe:

$$\nabla \Phi_{\alpha}(X) \approx \frac{1}{\alpha N} \sum_{\alpha N \text{ worst}} \frac{\nabla f_X(x_i)}{f_X(x_i)} x_i$$

- No!
- Why?

$$\nabla \Phi_{\alpha}(X) = \nabla \int_{-\infty}^{q} \alpha^{-1} f_{X}(x) x dx$$
$$= \int_{-\infty}^{q} \alpha^{-1} \nabla f_{X}(x) x dx + \alpha^{-1} \nabla q f_{X}(q) dx$$

Proposition

We have

$$\nabla \Phi_{\alpha}(R(X)) = \mathbb{E}\left[\left.\frac{\nabla f_{X}(X)}{f_{X}(X)}(R(X) - q)\right| R(X) \le q\right]$$

Estimation algorithm

$$\nabla \Phi_{\alpha}(R(X)) \approx \frac{1}{\alpha N} \sum_{\alpha N \text{ worst}} \frac{\nabla f_X(x_i)}{f_X(x_i)} (R(x_i) - \hat{q})$$

where \hat{q} is empirical quantile.

Guarantees

- Gradient estimate bias is $\mathcal{O}(N^{-1/2})$
- Convergence w.p. 1 of SGD to local CVaR optimum

Tetris

- Softmax policy, standard features (Thiery and Scherrer, 2009)
- Bonus for clearing multiple rows
- Compare standard policy gradient with CVaRSGD



Tetris

- Softmax policy, standard features (Thiery and Scherrer, 2009)
- Bonus for clearing multiple rows
- Compare standard policy gradient with CVaRSGD



Avg. reward: **451** vs. **414**

A. Tamar, UC Berkeley

5

<ロト < 回 > < 回 > < 回 > < 三 > 三 三

Tetris

- Softmax policy, standard features (Thiery and Scherrer, 2009)
- Bonus for clearing multiple rows
- Compare standard policy gradient with CVaRSGD



Avg. reward: **451** vs. **414** Reward CVaR: 323 vs. **394**

< □ > < □ > < □ > < □ > < □ > < □ >

Tetris

- Softmax policy, standard features (Thiery and Scherrer, 2009)
- Bonus for clearing multiple rows
- Compare standard policy gradient with CVaRSGD





Risk-Averse

Standard

Risk Averse RL - Summary

- Unified approach to safety
 - Unluckiness
 - Model mismatch
- Efficient algorithms
- Applications
 - Finance (e.g., options^[1])
 - Robotics (simulation to real-world^[2], safety)

[1] T., Mannor, Xu., ICML 2014[2] Rajeswaran et al., 2016

Part 2:

Value Iteration Networks

T., Wu, Thomas, Levine, Abbeel. NIPS 2016

Goal: autonomous robots





Solution: reinforcement learning?

Image credit: http://www.wellandgood.com/wp-content/uploads/2015/02/Shira-fridge.jpg

Deep RL success: visual input > action



Mnih V., et al. Nature 2015



And many more...

Levine S., et al. JMLR 2016

Deep RL success: visual input > action



Reactive policies:



A. Tamar, UC Berkeley

Generalization in RL

- Reactive policies can solve complex tasks
- But do they understand?
- A simple test: **generalization** on grid worlds


- Reactive policies can solve complex tasks
- But do they understand?
- A simple test: **generalization** on grid worlds



Train



Test



A. Tamar, UC Berkeley

Train



Test



Observation: reactive policies do not generalize well

A. Tamar, UC Berkeley

Why do reactive policies fail to generalize?

- A sequential task requires a planning computation
- RL gets around that learns a mapping



Q/advantage/expert: planning on training domains

Why do reactive policies fail to generalize?

- A sequential task requires a planning computation
- RL gets around that learns a mapping



- Q/advantage/expert: planning on training domains
- New task need a new plan

Value iteration networks:

- Learn to plan
- Generalize to unseen tasks

Model

Network that Can Plan



Challenges:

- 1. What to plan?
- 2. How to plan?

3. How to use the plan?4. How to learn?

A. Tamar, UC Berkeley

- Value function: state > long-term value
 - Shortest-distance to-go
 - Max. sum of future rewards from state



 $V(s) = \max_{\pi} \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^{t} r(s_{t}) | \pi, s_{0} = s\right]$

- Sufficient for reactive policy
 - Action that maximizes value = optimal
- How to calculate?

Value := $\max_{action} \mathbb{E} \{ reward + next value \}$

















Value Iteration = Convnet !

Value Iteration Network (VIN)



1. What to plan?

Value Iteration Network (VIN)



- 1. What to plan?
- 2. How to plan?

Value Iteration Network (VIN)



1. What to plan?

3. How to use the plan?

2. How to plan?



Backprop through the whole network!

Learn planning computation that's useful for predicting actions

A. Tamar, UC Berkeley

Results

Supervised learning



5000 random maps X 7 trajectories

Compare VINs with:

- Convnet inspired by DQN architecture
- Fully convolutional net (FCN)

(Mnih et al. Nature 2015)

(Long et al. CVPR 2015) A. Tamar, UC Berkeley

Prediction loss on test set



Success of reaching goal



A. Tamar, UC Berkeley







- Is it only 'depth'?
 - VIN with untied weights
 - Degrades performance, esp. with less data

Train using RL
TRPO (Schulman et al. *ICML* 2015)



- A VIN is just a neural network
- Compose with perception & control modules
- Natural image input: Mars navigation





- Continuous control
 - Discrete planning in continuous domain?
 - Let the network figure it out!



- Train: Guided policy search (unknown dynamics)
- Compare: conv-net

(Mnih et al. Nature 2015, Lillicrap et al. ICLR 2016)



⁽Levine & Abbeel, *NIPS 2014*)

- Train: Guided policy search (unknown dynamics)
- Compare: conv-net

(Mnih et al. Nature 2015, Lillicrap et al. ICLR 2016)





⁽Levine & Abbeel, *NIPS 2014*)

More Examples

Navigate website links to find query



(Nogueira & Cho, NIPS 2016)

Hierarchical composition



Plan on Approximate Graph

- Navigate website links to find query (Nogueira & Cho, NIPS 2016)
- Plan on approx. graph (1st+2nd level categories)
- Learn feature similarity mappings



Plan on Approximate Graph

- Navigate website links to find query (Nogueira & Cho, NIPS 2016)
- Plan on approx. graph (1st+2nd level categories)
- Learn feature similarity mappings



Plan on Approximate Graph

- Navigate website links to find query (Nogueira & Cho, NIPS 2016)
- Plan on approx. graph (1st+2nd level categories)
- Learn feature similarity mappings



Extensions

- Follow-up work (independent)
 - 1st person navigation



Extensions



Cognitive Mapping and Planning for Visual Navigation, S. Gupta, J. Davidson, S. Levine, R. Sukthankar, J. Malik, arXiv 2017

A. Tamar, UC Berkeley

Summary – Value Iteration Networks

- Learn to plan -> generalization
- VIN framework
 - Motivated by planning theory
 - Differentiable planner (VI = convnet)
 - Easy to compose with perception/control
- Challenges
 - Planning in high-dim spaces
Part 3:

Hindsight MPC

T., Thomas, Zhang, Levine, Abbeel. ICRA 2017

Motivation

- Robotic manipulation
- Accurate dynamics model not available
- MPC: well-established control method
- Recently: adaptive MPC for manipulation ^[1,2]

[1] Fu, Levine, Abbeel. IROS 2016[2] Lenz, Knepper, Saxena. RSS 2015



At time t:



Predict *H*-step trajectory:



Using current dynamics estimate & policy: $\hat{x}_{s+1} = \hat{f}^t(\hat{x}_s, \pi^t(\hat{x}_s))$

Linearize dynamics:



 $\hat{f}_s(x,u) \approx A_s x + B_s u$

Solve linear quadratic regulator (LQR):



Solve linear quadratic regulator (LQR):



Take action u_t



Update dynamics estimate Repeat...

Adaptive MPC

- Sub-optimality:
 - Prediction (dynamics) error
 - Finite horizon
- Typically H << T:</p>
 - Mitigate model errors
 - Computational feasibility

Question: how to **improve** in a repeated task?

Idea: the Hindsight Plan

Revisit time t (offline):



Revisit time t (offline):



Revisit time t (offline):



Adaptive MPC

Hindsight action \bar{u}_t



Adaptive MPC

Hindsight action \bar{u}_t



- MPC Sub-optimality:
 - ✓ Prediction (dynamics) error
 - ✓ Finite horizon
- Hindsight actions are better!
 - How to improve (online) MPC?

Change cost such that online MPC would mimic the hindsight actions!





Make the online MPC mimic the hindsight plan!



Make the online MPC mimic the hindsight plan!

Illustrative Example



Illustrative Example



Illustrative Example



Results – Peg Insertion



State space: 26D

- 7 DoF arm (7 angles + 7 ang. velocities) •
- 2 end-effector points (6 positions + 6 velocities) • **Dynamics: GMM**
- Prior free space dynamics •
- **Online adaptation** •

Results – PR2 Peg Insertion



Conclusion

- Decades of research on planning/control
- Control algorithm \rightarrow inductive bias



- Principled approach to network design
- Better generalization (change in task)
- New algorithmic ideas

Conclusion

Inductive bias



Conclusion

Inductive bias

(algorithm learning)



Thanks

- Collaborators
 - Pieter Abbeel (UCB), Yinlam Chow (Stanford), Mohammad Ghavamzadeh (Adobe), Yonatan Glassner (Technion), Sergey Levine (UCB), Shie Mannor (Technion), Marco Pavone (Stanford), Garrett Thomas (UCB), Yi Wu (UCB), Tianhao Zhang (UCB)
- Funding
 - Siemens
 - Viterbi Scholarship

Berkeley

Artificial Intelligence Research Laboratory