[IJRR'18]: [ICRA'17,RSS'17,CoRL'17]

3

Microsoft[®]

Research

Data Driven Planning via Imitation Learning

Sanjiban Choudhury¹

Mohak Bhardwaj², Sebastian Scherer², Ashish Kapoor³, Gireeja Ranade³ and Debadeepta Dey³

Carnegie Mellon THE ROBOTICS INSTITUTE



2

Research theme: Adaptation in planning

Robust navigation across diverse scenarios/platforms







Research theme: Adaptation in planning

Robust navigation across diverse scenarios/platforms Complex manipulation across diverse tasks



Why do we need adaptive planning?

Decades of research on planning algorithms with provable worst-case guarantees



What happens to the finite time performance of such algorithms on real-world problems?

Navigation in varying scenarios







Collier trophy finalist $[AHS'14, ICRA'15, '16, '17]_5$

Can't we use a single versatile planner?

Optimal planners, by ignoring context, are unable to succeed in real-time



Challenge: Good real-time performance

Real-time performance of a planner varies across scenarios (obstacle configs, connectivity of space, presence of poor local minima)

Given a finite time budget, planner cannot explore all configurations

Planning module needs to adapt to the distribution of problems the robot encounters

Two paradigms of adaptation

Black-box adaptation

(adaptation outside planner)

White-box adaptation

(adaptation inside planner)

Two paradigms of adaptation

Black-box adaptation

(adaptation outside planner)

E.g. Learn to select an ensemble of planners from a library [ICRA'15, '16]

Prior Work

Zucker et al. '09, Dey et al.'13, Jetchev et al'13, Pan et al.'14, Poffald et al.'14, Dragan et al.'11, Palmieri et al.'15, Morales et al.'05

White-box adaptation

(adaptation inside planner)

Long term autonomy with black-box adaptation



Single baseline planner fails 36% of the missions Adaptive planning enables continuous back-to-back trials (>2.5 hrs) No hand-tuning of parameters, No access to test UAV

Two paradigms of adaptation

Black-box adaptation

(adaptation outside planner)

E.g. Learn to select an ensemble of planners from a library [ICRA'15, '16]

Prior Work

Zucker et al. '09, Dey et al.'13, Jetchev et al'13, Pan et al.'14, Poffald et al.'14, Dragan et al.'11, Palmieri et al.'15, Morales et al.'05 White-box adaptation

(adaptation inside planner)

E.g. Learning a heuristic policy [CoRL'17]

Prior Work

Yoon et al.'06, Xu et al.'07, Arafee et al.'11, Thayer et al'11, Garrett et al'16., Aine et al.'16, Phillips et al'15, Paden et al.'16

E.g. Bayesian active edge evaluation of roadmaps [NIPS'17, IJCAI'18]

Prior Work

Burns et al.'05, Pan et al.'12, Choudhury et al.'16, Huh et al.'16, Javdani et al.'14, Dellin et al.'16, Narayan et al.'17

Two paradigms of adaptation

Black-box adaptation

(adaptation outside planner)

E.g. An adaptive meta-planner that selects a planner from a library to execute given present context [ICRA'15, '16]

Prior Work

Zucker et al. '09, Dey et al.'13, Jetchev et al'13, Pan et al.'14, Poffald et al.'14, Dragan et al.'11, Palmieri et al.'15, Morales et al.'05

White-box adaptation

(adaptation inside planner)

E.g. Learning a heuristic policy [CoRL'17]

Yoon et al.'06, X Garrett et al'16., 4 E.g.

evaluation in graph search [NIPS'17, IJCAI'18]

Prior Work

Burns et al.'05, Pan et al.'12, Choudhury et al.'16, Huh et al.'16, Javdani et al.'14, Dellin et al.'16, Narayan et al.'17

Are there other planning domains where performance varies with obstacle configurations?

Informative path planning

Plan an optimal sequence of sensing actions to accomplish task





Information gathering for mapping (Charrow et al.'15)

Information gathering for surface reconstruction (Isler et al.'16)



Microsoft[®] **Research** 14

Challenges: Adapt to different tasks

Configuration of objects determine the utility of a policy



Inspection of power lines (objects of interest concentrated in a line)



Inspection of tunnels (objects of interest distributed around robot)

Planning module should adapt strategy accordingly



of this talk









3. Open Questions

Heuristic policies in search based planning



Objective: Guide a search tree from start to goal to find a feasible path

Why does a heuristic need to be adaptive?

Prior work has mainly focussed on bounding solution quality by defining heuristics as distance estimate of cost-to-go (Pohl'70, Pearl'84)

Compute estimates using relaxation-based (Likhachev et al.'09, Dolgov et al.'08) or learning-based (Xu et al.'07, Garrett et al.'16, Aine et al.'15, Paden et al.'16) approaches

Problem: Small estimation error leads to excessive expansions (i.e. large planning times)



INFLATED EUCLIDEAN HEURISTIC

Heuristic gets trapped in 'bug trap' due to greediness LEARNT HEURISTIC POLICY





Worlds with 'bug traps'

Heuristic does not get trapped, searches along periphery

Problem: Minimize expansions in BFS



Problem: Minimize expansions in BFS



Compute a heuristic policy that maps search state to node to expand

 $\pi(\mathcal{O},\mathcal{C},\mathcal{I}) \to v \in \mathcal{O}$

 $\min \mathbb{E}_{\phi \sim P(\phi)} c(\pi, \phi)$ Objective: π Number of BFS Distribution over world maps expansions

Sequential decision making under uncertainty



Cast as a Partially Observable Markov Decision Process (POMDP)

State
$$s_t : \{\mathcal{O}_t, \phi\}$$
Cost $\begin{cases} 0 & \text{if } v_g \in \mathcal{O}_t \\ 1 & \text{otherwise} \end{cases}$ Observation $o_t : \{\mathcal{V}_{\text{succ}}, \mathcal{E}_{\text{inv}}\}$ Action $a_t : v \in \mathcal{O}_t$ Cost $\begin{cases} 1 & \text{otherwise} \end{cases}$ History $\psi_t : \{\mathcal{O}_t, \mathcal{C}_t, \mathcal{I}_t\}$

POMDP solvers such as POMCP (Silver et al.'10) or DESPOT (Somani et al.'13) require a lot of online effort - we want cheap policies!

Model-free approaches such as Q-learning (Watkins et al.'92) or REINFORCE (Williams et al.'92) are sample inefficient

Key Idea: Imitate a clairvoyant oracle



Oracle is "clairvoyant" as it can process the whole world map

A quick detour through imitation learning

(in the context of MDPs)

Open loop demonstrations Corrective feedback

Cost-to-go feedback

Open loop demonstrations Corrective feedback

Cost-to-go feedback

Behavior cloning: Do plain supervised learning!

$$\hat{\pi} = \arg\min_{\pi \in \Pi} \mathbb{E}_{s \sim P(s|\pi_{OR})} [\ell(\pi(s), \pi_{OR}(s))]$$
(expert dist.) (imitation loss)



Problem: Compounding errors due to train test mismatch!



Solution: DAGGER (Ross, Gordon and Bagnell, 11). Train policy with initial data, collect more data with learner, aggregate data, re-train and repeat!

Open loop demonstrations Corrective feedback

Cost-to-go feedback

 $\mathbb{E}_{s \sim P(s|\pi_{OR})}[\ell(\pi(s), \pi_{OR}(s))] \quad \mathbb{E}_{s \sim P(s|\pi)}[\ell(\pi(s), \pi_{OR}(s))]$ (expert dist.) (imitation loss)

(Behaviour cloning)

(learner dist.) (imitation loss)

(DAGGER)

DAGGER minimizes immediate imitation loss (0-1, MSE, etc) it does not reason about the cost-to-go of actions

E.g. What if there are many good actions? What if some actions are much worse than others?



Solution: AGGREVATE (Ross and Bagnell, 14). Same as DAGGER, but label is obtained from cost-to-go of expert roll-out

Open loop demonstrations Corrective feedback

(expert dist.) (imitation loss)

(Behaviour cloning)

 $\mathbb{E}_{s \sim P(s|\pi_{OR})}[\ell(\pi(s), \pi_{OR}(s))] \quad \mathbb{E}_{s \sim P(s|\pi)}[\ell(\pi(s), \pi_{OR}(s))] \quad \mathbb{E}_{s \sim P(s|\pi)}[Q^{OR}(s, \pi(s))]$

(learner dist.) (imitation loss)

(DAGGER)

Cost-to-go feedback

(learner dist.) (expert cost-to-go)

(AGGREVATE)

How do we use AGGREVATE for imitating a clairvoyant oracle?

SAIL: Search as Imitation Learning [CoRL'17]



 $\begin{array}{c} \operatorname{Sample \ a \ world \ \phi} \\ \operatorname{Sample \ a \ world \ \phi} \\ \operatorname{from \ database} P(\phi) \end{array}$

P	(ϕ)	

 $egin{aligned} & \operatorname{Roll-in}_{to \ ext{get}} \operatorname{history}_{\psi_t} \pi_{ ext{mix}} \ & \operatorname{mixture} \ of \ & \operatorname{learner} + \operatorname{oracle} \ & \operatorname{to} \ & \operatorname{get} \ & \operatorname{search} \ & \operatorname{state} \ & \psi_t : \{\mathcal{O}_t, \mathcal{C}_t, \mathcal{I}_t\} \end{aligned}$



Execute a random action a_t Execute random and featurize (ψ_t, a_t) as f_t action and extract features $\{\psi_t, a_t\} \rightarrow f_t$



Roll out with oracle τ_{OR} add { f_t , Q^{σ_R} } to data to get value $Q^{\pi_{OR}}$ Aggregate data { f_t , $Q^{\pi_{OR}}$ } Update learner

31

Problem: Why does this even work?



AGGREVATE performance bounds indicate large realizability error due to information mismatch - not helpful

Trick: Define a hallucinating oracle [IJRR'18]

$$\tilde{Q}_{T-t+1}^{\tilde{\pi}_{\mathrm{OR}}}(\psi,a) = \mathbb{E}_{s\sim P(s|\psi)} \left[Q_{T-t+1}^{\pi_{\mathrm{OR}}}(s,a) \right]$$

(Hallucinating oracle) (Current (Clairvoyant oracle) posterior)

Hallucinating oracle computes expect action value at each time-step

$$J(\hat{\pi}) \leq J(\tilde{\pi}_{OR}) + T(\epsilon_{class} + \epsilon_{reg}) + \mathcal{O}\left(\frac{T\log T}{N}\right) - T\epsilon_{or}$$
Learner Hallucinating Realizability error of
performance oracle hallucinating oracle
performance (which is small)
$$33$$

Simple datasets: SAIL outperforms baselines

			Learning			Non-learning approaches			
Dataset	Sample Worlds	SAIL	SL	CEM	\mathbf{QL}	$h_{\mathbf{EUC}}$	$h_{\mathbf{MAN}}$	A*	MHA*
Alternating Gaps		0.039	0.432	0.042	1.000	1.000	1.000	1.000	1.000
Single Bugtrap		0.158	0.214	0.057	1.000	0.184	0.192	1.000	0.286
Shifting Gaps		0.104	0.464	1.000	1.000	0.506	0.589	1.000	0.804
Forest		0.036	0.043	0.048	0.121	0.041	0.043	1.000	0.075
Bugtrap+Forest		0.147	0.384	0.182	1.000	0.410	0.337	1.000	0.467
Gaps+Forest		0.221	1.000	1.000	1.000	1.000	1.000	1.000	1.000
Mazes		0.103	0.238	0.479	0.399	0.185	0.171	1.000	0.279
Multiple Bugtraps		0.479	0.480	1.000	0.835	0.648	0.617	1.000	0.876

Table shows average normalized expansions

Evaluating SAIL on helicopter problems



An autonomous helicopter navigating in a canyon



Dataset of canyons

Learning a heuristic for 4D search (x,y,z,heading)



A* using dubins distance heuristic times out (2531 states, 7000ms)

SAIL expands 18 states in 100 ms

Applying SAIL on real world problems

High speed no-fly-zone avoidance requires real-time nonholonomic path planning



A* with Dubins heuristic times out



Expands 1910 states in time budget (1000 ms)

Dubins distance is a poor estimate of expansions-to-go

SAIL learns to follow maze wall

Flight test with onboard execution

(Choudhury, Bhardwaj, Maeta, Scherer)

of this talk

2. Informative path planning

3. Open Questions

What is budgeted information gathering?

World map is hidden and revealed through measurements

Exact solutions difficult to compute

Known world map

NP-Hard (Singh et al.'09, Arora et al.'17, Zhang et al.'16)

Hidden world map Belief over worlds conditioned on measurements (Lim et al.'16, Hollinger et al.'12, '17, Singh et al '09)

Myopic heuristics are surprisingly good!

Given history (of nodes visited and measurements received) Select action (node to visit)

Heuristics select node that maximize weighted information gain (Isler et al.'16)

Different heuristics use different weights (e.g Average entropy, Rear-side Voxel, Occlusion Aware)

Average Entropy heuristic is poor when objects concentrated in a cluster

Rear-Side Voxel heuristic is poor when objects distributed in world

Cast problem as a POMDP

State
$$s_t = (v_1, \dots, v_t, \phi)$$
vertices visited, world (hidden)Action $a_t = v_{t+1}$ next vertex to visitReward $R(s, a) = . \begin{cases} \Delta_{\mathcal{F}}(a \mid \xi, \phi) & \text{if } a \in \mathcal{A}_{\text{feas}}(s) \\ -\infty & \text{otherwise} \end{cases}$ utility gained

Observation $o_t = y_t$ measurement received

History $\psi_t = \{v_i, y_i\}_{i=1}^t$ history of actions, measurements

Key Idea: Imitate a clairvoyant oracle

Oracle: Generalized Cost Benefit (Zhang and Vorobeychik, 2016)

Imitate oracle that fully observes (clairvoyant) state and plans good actions

The ExpLOre algorithm [ICRA' 17]

Sample a world from the database $e^{\phi} \sim P(\phi)$

Roll-in a mixture of learner + oracle to get history $\psi_t = \{v_i, y_i\}_{i=1}^t$

Pick an actionRoll-out oracleExtract featureto get value $\{\psi_t, a_t\} \rightarrow f_t$ $Q^{\pi_{OR}}(s_t, a_t)$ Aggregate dataUpdate learner

Comparison across datasets / problems

IG heuristic baselines

Dataset	Sample World Maps	Problem	Learner	Average Entropy	Occlusion Aware	Unobserved Voxels	Rear Side Voxels	Rear Side Entropy
Concentrated Parallel Lines (2D)		HIDDEN-UNC HIDDEN-CON	(0.42, 0.45) (0.18, 0.27)	(0.20, 0.22) (0.16, 0.20)	(0.06, 0.09) (0.07, 0.09)	(0.20, 0.25) (0.14, 0.18)	(0.36, 0.41) (0.19, 0.24)	(0.30, 0.34) (0.21, 0.26)
Distributed Blocks (2D)		HIDDEN-UNC HIDDEN-CON	(0.37,0.41) (0.20, 0.26)	(0.26, 0.30) (0.21, 0.26)	(0.11, 0.16) (0.11, 0.16)	(0.22, 0.29) (0.15, 0.20)	(0.22, 0.29) (0.15, 0.18)	(0.24, 0.28) (0.16, 0.19)
Poisson Forest of Circular Discs (2D)		HIDDEN-UNC HIDDEN-CON	(0.58, 0.61) (0.54, 0.59)	$(0.59,0.62) \ (0.54,0.59)$	(0.49, 0.53) (0.42, 0.46)	(0.39, 0.45) (0.34, 0.41)	(0.53, 0.55) (0.37, 0.43)	(0.42, 0.47) (0.39, 0.44)
Tabular World of Rectilinear Blocks (2D)		Hidden-Unc Hidden-Con	$egin{array}{c} (0.43,0.53)\ (0.27,0.33) \end{array}$	(0.31, 0.35) (0.26, 0.29)	(0.20, 0.26) (0.18, 0.23)	(0.28, 0.35) (0.21, 0.28)	(0.35, 0.44) (0.18, 0.24)	(0.25, 0.31) (0.21, 0.27)
Bookshelves and Tables (3D)		Hidden-Unc Hidden-Con	$egin{array}{c} (0.14,0.31)\ (0.05,0.24) \end{array}$	(0.01, 0.04) (0.01, 0.04)	(0.01, 0.04) (0.01, 0.04)	(0.01, 0.04) (0.01, 0.04)	(0.01, 0.22) (0.01, 0.22)	(0.01, 0.19) (0.01, 0.19)
Cluttered Construction Site (3D)		Hidden-Unc Hidden-Con	$egin{array}{c} (0.14, 0.20) \ (0.08, 0.12) \end{array}$	(0.01, 0.12) (0.01, 0.12)	(0.01, 0.09) (0.01, 0.09)	(0.01, 0.09) (0.01, 0.09)	(0.01, 0.11) (0.01, 0.11)	(0.01, 0.10) (0.01, 0.10)
Office Desk and Chairs (3D)		HIDDEN-UNC HIDDEN-CON	$(0.69,0.80) \ (0.55,0.72)$	(0.46, 0.59) (0.46, 0.59)	(0.51, 0.63) (0.48, 0.63)	(0.51, 0.63) (0.48, 0.63)	(0.59, 0.67) (0.43, 0.52)	(0.61, 0.72) (0.41, 0.53)

ExpLOre consistent across environments

Special case when budget is infinity [RSS'17]

Simply imitating clairvoyant one-step-reward outperforms all baselines

Explanation: Can prove near-optimality due to adaptive submodularity (Golovin and Krause, '11)

Train on simulation, Test on real data

Train Data: Office desks created in Gazebo

Image: state state

Test Data: RGBD data (Sturm et al.)

of this talk

2. Informative path planning

3. Open Questions

Issues in data driven planning

Useless expansions!

Learnt policies are incredibly brittle!

We have reference policies that are not so brittle. Can we leverage this in a principled manner?

Where do we get data from?

Common use case: we have a few representative maps, difficult to figure out how to augment these maps to learn a non-brittle heuristic

Why can't we have universal policies?

We tried to train heuristics over all world maps and it led to degraded performance (History might not be sufficiently discriminatory to collapse posterior)