Meta-Learning Contextual Bandit Exploration

Amr Sharaf University of Maryland amr@cs.umd.edu

Hal Daumé III

Microsoft Research & University of Maryland me@hal3.name



Can we learn to explore in contextual bandits?









Training Mêlée by Imitation



Generalization: Meta-Features

- No direct dependency on the contexts x.
- Features include:
 - Calibrated predicted probability $p(a_t | f_t, x_t)$;
 - Entropy of the predicted probability distribution;
 - A one-hot encoding for the predicted action $ft(x_t)$;
 - Current time step t;
 - Average observed rewards for each action.



MÊLÉE	0	30	23	167	126	160	166	182		160
ε-greedy	-30	0	2	176	136	144	174	185		
ϵ -decreasing	-23	-2	0	177	136	141	176	184		80
EG ε-greedy	-167	-176	-177	0	-56	-3	57	48		0
LinUCB	-126	-136	-136	56	0	59	91	77		
τ-first	-160	-144	-141	3	-59	0	33	31		-80
Cover	-166	-174	-176	-57	-91	-33	0	-19		
Cover-nu	-182	-185	-184	-48	-77	-31	19	0		-160
	MÊLÉE	e-greedy	creasing	E-greedy	LinUCB	τ-first	Cover	Cover-nu		
Win			ε-de		S	S	St	a.	E	istics

Win statistics: each (row, column) entry shows the number of times the row algorithm won against the column, minus the number of losses.

MÊLÉE	0	30	23	167	1
ε-greedy	-30	0	2	176	1
ε -decreasing	-23	-2	0	177	1
EG ε-greedy	-167	-176	-177	0	-
				_	

Win / Loss Statistics

Win statistics: each (row, column) entry shows the number of times the row algorithm won against the column, minus the number of losses.

Theoretical Guarantees

- The no-regret property of Aggrevate can be leveraged in our meta-learning setting.
- We relate the regret of the learner to the overall regret of π .
- This shows that, if the underlying classifier improves sufficiently quickly, Mêlée will achieve sublinear regret.

Conclusion

- Q: Can we learn to explore in contextual bandits?
- A: Yes, by imitating an expert exploration policy;
- Generalize across bandit problems using meta-features;
- Outperform alternative strategies in most settings;
- We provide theoretical guarantees.