Caltech



California Institute of Technology

Co-Training for Policy Learning



Jialin Song



Yisong Yue



Hiro Ono

Policy Learning (Reinforcement & Imitation)



Policy Learning is Hard

- Long time horizons
- Sparse or expensive feedback
- Exponential in time horizon
 - (reinforcement learning)
- Infeasible to obtain sufficient demonstrations
 - (imitation learning)



Example: Learning to Search (Combinatorial Optimization)



[He et al., 2014] [Song et al., arXiv]

Learning from Multiple Views

Example: Minimum Vertex Cover



Graph View

[Khalil et al., 2017]

$$\max - \sum_{i=1}^{5} x_i,$$

subject to:

$$x_1 + x_2 \ge 1,$$

$$x_2 + x_3 \ge 1,$$

$$x_3 + x_4 \ge 1,$$

$$x_3 + x_5 \ge 1$$

 $x_4 + x_5 \ge 1,$

 $x_i \in \{0, 1\}, \forall i \in \{1, \cdots, 5\}$

Integer Program View (Branch & Bound View) [He et al., 2014]

Learning from Multiple Views

Example: Different Types of Integer Programs



ILP



QCQP

Co-Training [Blum & Mitchell, 1998]

- Many learning problems have different sources of information
- Webpage Classification: Words vs Hyperlinks



(Taken from Avrim Blum's slides)

Semi-Supervised Regression with Co-Training

Zhi-Hua Zhou and Ming Li
National Laboratory for Novel Software Technology
Nanjing University, Nanjing 210093, China
Minmin Chen, Kil

{zhouzh, lim}@lamda.nju.edu.cn

Co-Training for Domain Adaptation

Minmin Chen, Kilian Q. Weinberger Department of Computer Science and Engineering Washington University in St. Louis St. Louis, MO 63130 mc15, kilian@wustl.edu John C. Blitzer Google Research 1600 Amphitheatre Parkway Mountain View, CA 94043 blitzer@google.com Email Classification with Co-Training

Svetlana Kiritchenko and Stan Matwin

School of Information Technology and Engineering University of Ottawa Ottawa, ON, Canada {svkir,stan}@site.uottawa.ca

A New Analysis of Co-Training

			label	ed E Co-Trainin	ng and Expansion: Theory and Pra	ctice	dging	
Wei Wang Zhi-Hua Z National Ke	Vei Wang hi-Hua Z Sational Ke Applying Co-Training methods to Statistical Parsing* Anoop Sarkar Dept. of Computer and Information Science University of Pennsylvania 200 South 33rd Street, Philadelphia, PA 19104-6389 USA anoord line cis upper edu Bayesian Co-Training		- DU.CN EDU.CN	Maria-Florina	Understanding the Behavi		ior of Co-training	
			Computer Scien Carnegie Mello Applying Co-Training to Referen		Kamal Nigam		Rayid Ghani ^ :hool of Computer Science arnegie Mellon University Pittsburgh, PA 15213 rayid@cs.cmu.edu	
			ristoph Müller	Sony International (Europe) Gombur Advanced Technology Center S Heinrich-Hertz-Straße 1 70327 Stuttgart, Germany	Michael Strube European Media Laboratory CombH Unsupervised Improvement of Visual Detect			
			Villa Bosch				ectors using Co-Training	
Shipeng Yu Balaji Krishnapuram Business Intelligence and Analytics Siemens Medical Solutions USA, Inc. Si Valley Stream Parkway		SHIPENG.YU@SIEMENS.COM BALAJI.KRISHNAPURAM@SIEMENS.COM	Volfsbrunnenweg 33 eidelberg, Germany ml.villa-bosch.de		Anat Levin* School of CS and Eng. The Hebrew University 91904 Jerusalem, Israel alevin@cs.huji.ac.il	Paul Viola Microsoft Research One Microsoft Way Redmond, WA 98052 viola@microsoft.com	Yoav Freund Computer Science Dept. Columbia University New York, NY 10027 freund@cs.columbia.edu	
Malvern, PA	1 19355, USÁ			Reinforced Co-Training				
Rómer Ros ^Y ⁴ _S R B 	A Co-training Approach for Multi-view Spectral Clus		ering	PAC Generalization Bo	unds for Co-training	- Depar T Sant	William Yang Wang rtment of Computer Science Jniversity of California a Barbara, CA 93106 USA	
Si 5. M Ab Hai Dep	bhishek Kumar Il Daumé III partment of Computer Sci	ABHISHEK© HAL@UMI ience, University of Maryland, College Park, MD 20742, USA	CS.UMD.EDU ACS.UMD.EDU	Sanjoy Dasgupta Michael L AT&T Labs-Research AT&T Labs dasgupta@research.att.com mlittman@res	. Littman David McAllester s-Research AT&T Labs-Research earch.att.com dmac@research.att.com	— mwi	lliam@cs.ucsb.edu	

Semi-Supervised Regression with Co-Training

Zhi-Hua Zhou and Ming Li National Laboratory for Novel Software Technology Nanjing University, Nanjing 210093, China {zhouzh, lim}@lamda.nju.edu.cn

Co-Training for Domain Adaptation

Minmin Chen, Kilian Q. Weinberger Department of Computer Science and Engineering Washington University in St. Louis John C. Blitzer Google Research 1600 Amphitheatre Parkway

Email Classification with Co-Training

Svetlana Kiritchenko and Stan Matwin

School of Information Technology and Engineering University of Ottawa Ottawa, ON, Canada {svkir,stan}@site.uottawa.ca

g

Wei Wang Zhi-Hua Z National Ke

Co-training for Policy Learning

Shipe	Jialin Song [†]	Ravi Lanka[‡] † California Inst	Yisong Yue [†] itute of Technology	Masahiro Ono [‡]	ining
Balaj Busino Sieme 51 Val Malve	K ss I s A ey. n,	lsion Laboratory, C	California Institute of	f Technology	u
Róm Yi 4	r Kosiekkue inito	-111			William Yang Wang
R B Si 5. M Hal Depa	A Co-training Approach for Multi-view Spec	PAC Genera	lization Bounds for Co-training	Department of Computer Science University of California Santa Barbara, CA 93106 USA	
	Abhishek Kumar Hal Daumé III Department of Computer Science, University of Maryland, College Park, MD 20	ABHISHEK@CS.UMD.EDU HAL@UMIACS.UMD.EDU 1742, USA	Sanjoy Dasgupta AT&T Labs-Research dasgupta@research.att.co	Michael L. Littman David McAllester AT&T Labs–Research AT&T Labs–Research m mlittman@research.att.com dmac@research.att.com	

What's Different about Policy Co-Training?

Sequential Decisions vs 1-Shot Decisions



[1] "Learning combinatorial optimization algorithms over graphs" [Khalil et al., 2017]
[2] "Learning to Search in Branch and Bound Algorithms" [He et al., 2014]
[3] "Learning to Search via Retrospective Imitation" [Song et al., 2019]

Intuition



[1] "Learning combinatorial optimization algorithms over graphs" [Khalil et al., 2017]
[2] "Learning to Search in Branch and Bound Algorithms" [He et al., 2014]
[3] "Learning to Search via Retrospective Imitation" [Song et al., 2019]

Intuition



[1] "Learning combinatorial optimization algorithms over graphs" [Khalil et al., 2017] [2] "Learning to Search in Branch and Bound Algorithms" [He et al., 2014] Intuition [3] "Learning to Search via Retrospective Imitation" [Song et al., 2019] π^1 (5)225 E.g., [1] **MVC** Instance Demonstration $\max - \sum_{i=1}^{n} x_i,$ $x_1=0$ subject to: E.g., [2,3] π^2 $x_2 = 1$ $x_1 + x_2 \ge 1,$ **Better!** $x_2 + x_3 \ge 1,$ $x_3 = 1$ $x_3 + x_4 \ge 1,$ $x_4=1$ $x_5=0$ $x_3 + x_5 \ge 1,$ $x_4 + x_5 \ge 1,$ $x_i \in \{0, 1\}, \forall i \in \{1, \cdots, 5\}$

Theoretical Insight

- Different representations differ in hardness
- Goal: quantify improvement



(Towards) a Theory of Policy Co-Training

- Two MDP "views": $M^1 \& M^2$ • $f^{1 \to 2}(\tau^1) \Longrightarrow \tau^2$ (and vice versa)
 - Realizing τ^1 on $M^1 \Leftrightarrow$ realizing τ^2 on M^2



- Question: when does having two views/policies help?
 - Policy Improvement (next slide)
 - Builds upon [Kang et al., ICML 2018]
 - Optimality Gap for Shared Action Spaces (in paper)
 - Builds upon [DasGupta et al., NeurIPS 2002]



Builds upon theoretical results from [Kang et al., ICML 2018]

Policy Improvement Bound (Summary)

$$J(\pi^{\prime 1}) \ge J_{\pi^1}(\pi^{\prime 1}) - \frac{2\gamma \left(\alpha_{\Omega}^1 \varepsilon_{\Omega}^1 + 4\beta_{\Omega_2}^2 \varepsilon_{\Omega_2}^2\right)}{(1-\gamma)^2} + \delta_{\Omega_2}^2$$

- Minimizing $\beta_{\Omega_2}^2 \rightarrow \text{low disagreement between } \pi^2 \text{ vs } \pi^1$
- Maximizing $\delta_{\Omega_2}^2 \rightarrow$ high performance gap π^2 over π^1 on some MDPs

CoPiEr Algorithm (Co-training for Policy Learning)



Performance comparison for Minimum Vertex Cover



Co-Training for Policy Learning (summary)

- First formal framework for policy co-training
- Novel theoretical insights
- Principled algorithm design
- Strong experimental results

$$J(\pi^{\prime 1}) \geq J_{\pi^1}(\pi^{\prime 1}) - \frac{2\gamma \left(\alpha_{\Omega}^1 \varepsilon_{\Omega}^1 + 4\beta_{\Omega_2}^2 \varepsilon_{\Omega_2}^2\right)}{(1-\gamma)^2} + \delta_{\Omega_2}^2$$



References

- "Co-Training for Policy Learning," Jialin Song, Ravi Lanka, Yisong Yue, Masahiro Ono, UAI 2019
- "Combining Labeled and Unlabeled data with Co-training," Avrim Blum, Tom Mitchell, COLT 1998
- "PAC Generalization Bounds for Co-training," Sanjoy DasGupta, Michael Littman, David McAllester, NeurIPS 2002
- "Policy Optimization with Demonstrations," Bingyi Kang, Zequn Jie, Jiashi Feng, ICML 2018
- "Learning Combinatorial Optimization over Graphs," Elias Khalil, Hanjun Dai, Yuyu Zhang, Bistra Dilkina, Le Song, NeurIPS 2017
- "Learning to Search in Branch and Bound Algorithms," He He, Hal Daume III, Jason Eisner, NeurIPS 2014
- "Learning to Search via Retrospective Imitation," Jialin Song, Ravi Lanka, Albert Zhao, Aadyot Bhatnagar, Yisong Yue, Masahiro Ono, arXiv

Extra Slides



Policy Improvement Bound (detailed)



Builds upon theoretical results from [Kang et al., ICML 2018]

Performance comparison for Risk-Aware Path Planning





OpenAl Gym & Mujoco

Partitioned state space into two views

Shared action space

RL on both views