Near-Optimal Machine Teaching via Explanatory Teaching Sets

Yuxin Chen Oisin Mac Aodha Shihan Su Pietro Perona

Yisong Yue

Caltech



Our Approach

Explanation-based Machine Teaching

Baird Sparrow Chipping Sparrow Chipping Sparrow Baird Sparrow

Surrogate objective function:

edge weight after observing S

remaining weight of the bipartite graph upon observing set S:

 $w(\{f,g\} \mid \mathcal{S}) = w(\{f,g\}) \prod \mathbb{P}(\mathbf{e}_i, y_i \mid f, g)$

Bipartite graph constructed by NOTES. Size of nodes represents prior prob. Edges are drawn between f^* (resp. g^*) and all g's (resp. all f's).

Near-Optimal Teaching via Explanatory Sets (NOTES)

Input Teaching image set $\{(x, y, e)\}_{1:m}$; hyp. $\{h(x):=g \circ f(x)\}$;



Student's ability to learn a new concept can be greatly improved by providing them with clear and interpretable **explanations** from a knowledgeable teacher

 $r(\mathcal{S}) = \sum_{\{f,g\} \in \mathcal{E}} w(\{f,g\} \mid \mathcal{S})$

Thm The worst-case cost of NOTES achieving error ϵ is within a logarithmic factor of the worstcase cost of the optimal algorithm achieving error of at least $P(f^*)P(g^*)\epsilon/2$

-	noise params $\{v^{f}, v^{g}\}_{I:m}$; prior $P(f), P(g)$; tolerance ϵ	
Output	Selected images to teach, S	

Start	$S \leftarrow \varnothing;$
Loop	$i^* = \operatorname{argmin}_i r(S \cup \{i\});$
	S ← S ∪ {i};
Until	$r(S) \leq P(f^*)P(g^*) \epsilon$

Experimental Results

Datasets



Objects from Mars have blue top, square middle and thick base Objects with any other combination belong to Jupiter

Vespula vs Weevil



Weevil has a mid-sized body and head, while Vespula does not.

Caltech-UCSD Birds (CUB) dataset

- CUB-1: Western Gull vs Heermann Gull
- CUB-2: Baird Sparrow vs Chipping Sparrow
- CUB-3: Song Sparrow vs Northern Waterthrush



black leg

grey bell





white breast buff beak striped breast pattern black beak white belly vellow beak

solid breast pattern striped wing striped back brown undertail short beak striped tai

solid wing

solid tail

all-purpose beak

Explaining Teaching Examples

LIME Explainer (Ribeiro et al., 2016)

LIME aims to reproduce the predictive results of h* in the vicinity of the input via a sparse linear classifier



"X" is classified as the blue class, because its horizon axis is less than the threshold marked by the red line

Likelihood functions for inconsistent explanations



Results

simulated and real-human learners



Baselines:

random: Rand. examples with no explanations randexp: Rand. examples with rand. explanations STRICT: Label-based greedy approach

With explanation-based machine teaching, learners achieve

- Better accuracy
- Faster question answering at test time