LABEL-FREE CONCEPT BOTTLENECK MODELS

ICLR 2023

XAI in Computer Vision

Feature Attribution vs Concept-based XAI



XAI in Computer Vision

Landscape of Concept-based XAI model in Computer vision

Base	Model	Structure	Additional training	Concept Supervision	Concept Granularity	Concept Form	Model Agnostic	Backbone
	Network Dissection (2017)	Post-hoc	No	Weakly-supervised	Neuron	Explicit	No	CNN
Dissection	GAN-Dissect (2019)	Post-hoc	No	Unsupervised	Neuron	Explicit	No	GAN (generator)
	CLIP-Dissect (2022)	Post-hoc	No	Weakly-supervised	Neuron	Explicit	No	CLIP
CAV	TCAV (2018)	Post-hoc	Yes	Fully-supervised	Image	Semi-Explicit	Yes	CNN
ACE	ACE (2019)	Post-hoc	Yes	Unsupervised	Image	Implicit	No	CNN
	CRAFT (2023)	Post-hoc	Yes	Unsupervised	Image	Implicit	Yes	CNN
	CBM (2020)	Pre-	hoc	Fully-supervised	Image	Explicit	No	CNN
	Post-hoc CBM + CAV (2022)	Post-hoc	Yes	Fully-supervised	Image	Explicit	Yes	CNN
	Post-hoc CBM + CLIP (2022)	Post-hoc	No	Weakly-supervised	Image	Explicit	No	CLIP
	Label-Free CBM (2023)	Pre-	-hoc	Weakly-supervised	Image	Explicit	Yes	CNN

Label-Free CBM

CBM vs Label-Free CBM



Label-Free CBM

CBM vs Label-Free CBM



CHAPTER 2 Label-Free CBM

Architecture



Figure 2: Overview of our pipeline for creating label-free CBM.

Label-Free CBM



Figure 2: Overview of our pipeline for creating label-free CBM.

Experimental setup

Baseline				
Duschine	Model	Concept Usage	Concept Supervision	Structure
	СВМ	0	Fully-supervised	Pre-hoc
	Post-hoc CBM (P-CBM)	0	Fully-supervised	Post-hoc
	Post-hoc CBM (w/ Concept Net)	0	Weakly-supervised	Post-hoc
	Sparse Linear Model	Х	_	Post-hoc
	IBD	Х	_	Post-hoc
	Label-Free CBM	0	Weakly-supervised	Pre-hoc
• CIFAR-10 • CIFAR-100 • CUB-200 • Places365 • ImageNet		(1) A (2) G (3) In (4) M	ccuracy lobal Explainabil stance-level Exp anual Editing of	ity Iainability Final Laye

Result

	(I) Flexibility		(II) Inter	rpretability	(III) Performance	
Method:	Without labeled concept data	Any network architecture	Sparse final layer	All features interpretable	Preserves accuracy	Extends to ImageNet scale
CBM	No	Yes	No	Yes	No	No
IBD	No	Yes	No	No	Yes	No
P-CBM	No	Yes	Yes	Yes	No	No
P-CBM (CLIP)	Yes	No	Yes	Yes	No	Maybe
P-CBM-h	No	Yes	Yes	No	Yes	No
P-CBM-h (CLIP)	Yes	No	Yes	No	Yes	Maybe
Label-free CBM (This work)	Yes	Yes	Yes	Yes	Yes	Yes

				Dataset		
Model	Sparse final layer	CIFAR10	CIFAR100	CUB200	Places365	ImageNet
Standard	No	88.80%*	70.10%*	76.70%	48.56%	76.13%
Standard (sparse) P-CBM P-CBM (CLIP)	Yes Yes Yes	82.96% 70.50%* 84.50%*	58.34% 43.20%* 56.00%*	75.96% 59.60%* N/A	38.46% N/A N/A	74.35% N/A N/A
Label-free CBM (Ours)	Yes	86.40% ± 0.06%	65.13% ± 0.12%	$74.31\% \\ \pm 0.29\%$	43.68% ± 0.10%	$71.95\% \pm 0.05\%$

Result

ImageNet CBM Orange vs Lemon



Places365 CBM Mountain vs Mountain Snowy

Concept	Prediction
 arid climate rocky and dry a large, sheer rockface 	
a crater	
- has a crater at the top	
a deep, narrow valley	
a gorge	mountain
a hiking boot	
a cone-shaped mountain	
a high elevation	
a large, rocky peak	
snow-covered slopes	mountain snowy
a large, flat expanse of snow	
a large, flat piece of ice	
a ski patrol	
ice	
may have snow or ice on top	
a summit	
minarets	

Result



Pred:Red headed Woodpecker - Conf: 0.706 - Logit:7.80 - Bias:-0.12



Pred:junkyard - Conf: 0.417 - Logit:10.09 - Bias:-0.30 +2.56 junk +2.05 old cars +1.11 debris +0.75 car parts a excavator +0.5a windshield +0.48filled with trash or debris +0.4Sum of 2200 other features +2.56 1.0 1.5 2.0 Concept contributions 0.0 0.5 2.5



Places 365

CUB