Information Visualization

eXplainable Artificial Intelligence

Dr. David Koop







D. Koop, CSCI 628, Fall 2021







Feature Visualization



Edges (layer conv2d0)

Textures (layer mixed3a)

D. Koop, CSCI 628, Fall 2021









Objects (layers mixed4d & mixed4e)





Northern Illinois University







Feature Visualization vs. Attribution



Feature visualization answers questions about what a network—or parts of a network—are looking for by generating examples.

D. Koop, CSCI 628, Fall 2021





Attribution ¹ studies what part of an example is responsible for the network activating a particular way.

[<u>C. Olah et al.</u>, 2017]









Feature Vis by Optimization

- "[W]hat kind of input would cause a certain behavior"
- Start from random noise and iteratively tweak (using derivatives)



• What are the objectives? (Where are we going?) - Neuron, channel, layer (has DeepDream "interesting" objective

D. Koop, CSCI 628, Fall 2021



Step 48













Optimization Objectives

Different **optimization** objectives show what different parts of a network are looking for.

- **n** layer index
- x,y spatial position
- **z** channel index
- **k** class index



Neuron layer_n[x,y,z]



Channel layer_n[:,:,z]





D. Koop, CSCI 628, Fall 2021









Layer/DeepDream layer_n[:,:,:]²



Class Logits pre_softmax[k]



Class Probability softmax[k]











Why not Examples?

Dataset Examples show us what neurons respond to in practice

Optimization isolates the causes of behavior from mere correlations. A

neuron may not be detecting what you initially thought.









Animal faces—or snouts? mixed4a, Unit 240

D. Koop, CSCI 628, Fall 2021

Clouds—or fluffiness? mixed4a, Unit 453

Buildings—or sky? mixed4a, Unit 492







Diversity

- Examples can be diverse
- Optimization may give us one very positive takeaway
- Add a diversity term!



Simple Optimization



Optimization with diversity reveals multiple types of balls. Layer mixed5a, Unit 9

D. Koop, CSCI 628, Fall 2021



Dataset examples







Naive Optimization

Even if you carefully tune learning rate, you'll get noise.

Optimization results are enlarged to show detail and artifacts.



REPRODUCE IN A CO NOTEBOOK



Step 1

Step 32

D. Koop, CSCI 628, Fall 2021

Step 128

Step 256

Step 2048











Regularization to Avoid Noise

- Frequency penalization: penalize high-frequency noise Transformation robustness: jitter/rotate/scale images and still activate
- Learned priors: learn a prior and try to enforce it







The Building Blocks of Interpretability

The Building Blocks of Interpretability

Interpretability techniques are normally studied in isolation. We explore the powerful interfaces that arise when you combine them and the rich structure of this combinatorial space.

CHOOSE AN INPUT IMAGE

For instance, by combining feature visualization (*what is a neuron looking for?*) with attribution (*how does it affect the output?*), we can explore how the network decides between labels like **Labrador retriever** and **tiger cat**.

CHANNELS THAT MOST SUPPORT ...

LABRADOR RETRIEVER

feature visualization of channel hover for attribution maps →			
net evidence	1.63	1.51	
for "Labrador retriever"	1.22	1.24	
for "tiger cat"	-0.40	-0.27	

D. Koop, CSCI 628, Fall 2021



<u>Several floppy ear</u> detectors seem to be important when distinguishing dogs, whereas <u>pointy ears</u> are used to classify "tiger cat".



1.32 0.13



TIGER CAT





Combine Feature Vis and Activation



Activation Vector

Channels







Feature Vis + Attribution













Spatial Attribution with Saliency Maps

INPUT IMAGE



OUTPUT CLASSES

_abrador Retriever	
Golden Retriever	
Tennis Ball	1
Rhodesian Ridgeback	
Appenzeller	

mixed3a



mixed4a



D. Koop, CSCI 628, Fall 2021

OUTPUT FACTORS

Labrador Retriever	Tiger	
Golden Retriever	Tiger Cat	
Beagle	Lynx	
Kuvasz	Collie	
Redbone	Border Collie	

mixed4d

REAL OF	parts.	100	100	52	1285	CEN	1000	2657	1623	100	153	F7306	NOR .
	5.000		Sec.	1		K3	83		250		398	125	10
		0				8		128	85	1	20	*	
								5			251		
												5	
		** 1			8								
			10	12	192	1		R		2	<u>(</u>	51	
					20			21				M	
							M	5				5	
		2					2	3 2					
20	8				63	EZ			5			8 2	
	8												
					-		105						

mixed5a











Channel Attribution

INPUT IMAGE



OUTPUT CLASSES

Labrador Retriever	
Golden Retriever	
Tennis Ball	1
Rhodesian Ridge	
Appenzeller	

TOP CHANNELS SUPPORTING LABRADOR RETRIEVER













Showing 3 of **512**





Showing 3 of 480

• • •

Showing 3 of 508

D. Koop, CSCI 628, Fall 2021

[C. Olah et al.]



Showing 3 of 512









Factoring into Neuron Groups

INPUT IMAGE

By using non-negative matrix factorization we can reduce the large number of neurons to a small set of groups that concisely summarize the story of the network.





NEURON GROUPS based on matrix factorization of mixed4d layer



EFFECT of neuron groups on output classes

Labrador retriever	2.249	3.755	-1.193	-1.141	1.117	-1.892
beagle	3.298	0.599	-0.110	-0.356	-0.133	-2.618
tiger cat	-0.350	-0.994	-1.607	0.116	0.248	0.205
lynx	0.111	-0.642	-0.057	0.117	1.120	0.152
tennis ball	0.920	1.336	0.152	-0.885	1.227	-0.480

D. Koop, CSCI 628, Fall 2021

ACTIVATIONS of neuron groups



6 groups









Adding InfoVis

INPUT IMAGE



To understand multiple layers together, we would like each layer's factorization to be "compatible"—to have the groups of earlier layers naturally compose into the groups of later layers. This is also something we can optimize the factorization for.

- positive influence
- negative influence



D. Koop, CSCI 628, Fall 2021









Progress Reports



