Reliable and Interpretable Artificial Intelligence

Lecture 10: Visualization

Marc Fischer

ETH Zurich

Fall 2020



http://www.sri.inf.ethz.ch



Eidgenössische Technische Hochschule Zürich Swiss Federal Institute of Technology Zurich

Aspects of Reliable AI



AI Safety & Robustness

Verification of safety and robustness properties; training that enforces them.

e.g. Adversarial Robustness, Robust Training, Querying and Constraining Neural Nets (covered in upcoming and previous lectures)



Interpretability

Understanding the models decision making process.

e.g. Attention Mechanisms, Neurosymbolic Approaches, Visualization of decision making (today)



Privacy & Fairness

Data of individuals in produced and they are not discriminated,

e.g. Differential privacy, Enforcing Fairness in Neural Networks (covered in upcoming lectures)

Today: Visualizing Classification

INPUT IMAGE



ACTIVATIONS of neuron groups

Convolutional Neural Net (CNN)



- Convolutional layers: $f^{(k)}(x) = ReLU(x * F + b)$
- Linear layers: $f^{(k)}(x) = ReLU(A \cdot x + b)$

Example: 2D Convolution

0	0	0.8	0				
					0	0	0
0	1.0	0	0		1.0	0	1.0
0	0.3	0	1.0	*	1.0	0	1.0
-		-			0	0	0
0	0	0	0				
0	0	0	0		0	0	0

 $x \in \mathbb{R}^{4 \times 4}$

 $F \in \mathbb{R}^{3 \times 3}$

Example: 2D Convolution



 $x \in \mathbb{R}^{4 \times 4}$

 $F \in \mathbb{R}^{3 \times 3}$

Example: 2D Convolution

0 0

3

3 -4



Output Volume (3x3x2) Filter W1 (3x3x3) w1[:,:,0] 0[:,:,0] 0 0 0 1 -2 3 2 -1 0 0 -1 -3 0 w1[:,:,1] 0[:,:,1] 0 1 0 -1 -1 1 1 2 1 w1[:,:,2] -1 1 1

When we apply a convolutional layer we in fact apply many convolutions in parallel

Bias b1 (1x1x1) b1[:,:,0] 0

-1 0 -1

0 1 -1

We refer to the to the different results as channels. The above green output has two channels (as two filters were used).

toggle movement

https://cs231n.github.io/convolutional-networks/

Convolution

- Convolutions (formally cross-correlations) have a high response if the local pattern of the filter (kernel) matches the pattern in the data
- For example: a matrix encoding high contrast between left and right (shown below) yields a vertical edge detector



$$\begin{bmatrix} -1 & 0 & +1 \\ -2 & 0 & +2 \\ -1 & 0 & +1 \end{bmatrix}$$



AlexNet

- 2012 Krizhevsky et al. "AlexNet"
 - low error on ImageNet benchmark
 - Observation: feature similarity

"If two images produce feature activation vectors with a **small Euclidean separation**, we can say that the higher levels of the neural network consider them to be **similar**."

AlexNet

- Let R(x) be the vector of neuron activations for a layer deep in the network
- If the distance $||R(x_1) R(x_2)||_2$ is small, then the images x_1, x_2 are similar

1st column: image from test set

others: 6 training set images with most similar R(x)



Features

- Consequence: neuron activations capture semantic features of the input
- Feature Visualization: For a given neural network, what are these features?
- Feature Attribution: Which areas of a given image are responsible for classification in a given neural network?

AlexNet: First Layer Features

- Weights of the first convolutional layer (11x11x3) shown as color images
- Detectors for edges, simple patterns and color blobs



Inverting Convolutions

- Inspect neural networks in order to improve them
- For a given input and convolutional filter:
 - inverts convolutions to showcase the structure in the input causing high activation
 - first propagate input forward to the filter
 - then propagate these activations back to the input layer in the inverted network (shown next)
 - the outputs highlight which part of the image are responsible for the activation of the filter

Inverting Convolutions



Inverting Convolutions

- For a random set of convolutional filters in each layer, this shows the inputs causing the highest activations in the test set
- Pattern on the left reveals structure causing activation









find maximize $\boldsymbol{\chi}$ $score(x) - \lambda_1 R_1(x) - \dots - \lambda_m R_m(x)$

where

$$score(x) = mean(layer_n[x, y, z])$$

Find input *x* that maximizes score and minimized regularizes

Score denotes the activation of a neuron, channel or layer

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_[x,y,z]

Channel layer_[:,:,z]

Layer/DeepDream $layer_{1}[:,:,:]^{2}$



pre softmax[k]

Class Probability softmax[k]

 \rightarrow



Step 0



Step 4

 \rightarrow



Step 48



Step 2048

 \rightarrow



Olah et al., 2017

- Regularizers are crucial
 - Encode the prior "input should look a real image"
 - Optimization in Fourier basis, data whitening, added jitter etc.
- Optimize for multiple inputs at once
 - add diversity term loss to discourage similarity
 - multiple visualizations for one neuron



Simple Optimization



Optimization with diversity reveals four different, curvy facets. Layer mixed4a, Unit 97



Edges (layer conv2d0)

Textures (layer mixed3a)

Patterns (layer mixed4a)

Parts (layers mixed4b & mixed4c)

Objects (layers mixed4d & mixed4e)



Baseball—or stripes? *mixed4a, Unit 6*

Animal faces—or snouts? *mixed4a, Unit 240*

Clouds—or fluffiness? mixed4a, Unit 453

Buildings—or sky? mixed4a, Unit 492

Feature Attribution

- Which areas of the image are responsible for classification?
- Many techniques
 - Gradient based
 - Shapley values
 - Occlusion based
 - Etc.



Original Image

Grad-CAM 'Cat'



Original Image

Grad-CAM 'Dog'

Feature Attribution: Gradient Based

Junco Bird

Corn

Wheaten Terrier



Gradient of target logit w.r.t. the input shows how changing single pixels influences the logit

 $\partial logit_t(x)$ ∂x

Feature Attribution: Gradient Based



- To obtain a clearer result
 - Different methods apply scaling/conditioning of the gradient
- Just looking at the visualization is dangerous: many of these methods were found not to depend (much) on the actual weights in higher layers of the network
 - the gradient, by it's definition, however does

Feature Attribution: Shapley Values (general)

- Input x has a set of features P (e.g. individual pixels)
- Want to know how much each feature contributes to a function f(P) (e.g. a class logit)
- Classic result from game theory
- Theoretically justified



Feature Attribution: Shapley Values (image classification)

- Treat each pixel location as a feature
- Instantiate f with the logit of the target class (or the classification loss)
- Define meaning of f(S)
 - E.g. pixels locations not contained in S are set to a baseline value (zero, mean pixel value etc.)
- Computationally expensive
 - Exponentially many subsets
 - Can be approximated in polynomial time



Example: Shapley Values

- Want to compute the influence of x_1, x_2 on l_1 for $x = \binom{1}{1}$, R denotes ReLU
- Define f(S) as the value of l₁ when inputs not in S are set to 0



•
$$P = \{x_1, x_2\}$$

 $C_i = \sum_{S \subseteq P \setminus \{i\}} \frac{|S|! (|P| - |S| - 1)!}{|P|!} [f(S \cup \{i\}) - f(S)]$
 $\emptyset = 0$
 $\kappa(S)$

$$\{x_1\}$$
 1

D = (a a b)

$$\begin{cases} x_2 \} & 1 \\ x_{1,} x_2 \} \end{cases} \begin{array}{l} C_{x_1} = \kappa(\emptyset) \left[f(\{x_1\}) - f(\emptyset) \right] + \kappa(\{x_2\}) \left[f(\{x_{1,}, x_2\}) - f(x_2) \right] \\ = \frac{0!(2 - 0 - 1)!}{2!} [1 - 0] + \frac{1!(2 - 1 - 1)!}{2!} [2.5 - 1] = \frac{1}{2} (1 + 1.5) = 1.25 \\ C_{x_2} = 1.25 \end{array}$$

 x_1, x_2 both contribute 1.25 to $l_1 = 2.5$ for $x = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$

Combining Viualization & Attribution

- 1. Compute activations in a convolutional layer
- 2 Clusters activations k groups (here k = 6)
- 3. For each group:
 - show which pixel affect 1. it most by upscaling the activation map (grouped activations)
 - 2. show a visualization for the group of features

This allows for "humanscale" explanations (k groups rather than thousands)

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

REPRODUCE IN A CO NOTEBOOK

color key

INPUT IMAGE



ACTIVATIONS of neuron groups



NEURON GROUPS based on matrix factorization of mixed4d layer

6 groups



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

Combining Viualization & Attribution

- 1. Compute activations in a convolutional layer
- Clusters activations k groups (here k = 6)
- 3. For each group:
 - show which pixel affect it most by upscaling the activation map (grouped activations)
 - 2. show a visualization for the group of features

This allows for "humanscale" explanations (k groups rather than thousands)



How Do Adversarial Examples Impact Visualization?

Problem: Adversarial Examples

- Recall: If the distance $||R(x_1) R(x_2)||_2$ is small, then the images x_1, x_2 are similar
- Yet x and it's adversarial example x' are similar, but R(x) and R(x') are different
- Can we even find dissimilar x_1, x_2 with similar $R(x_1)$ and $R(x_2)$?

Counterexample

find
$$x'_1 \coloneqq x_1 + \delta$$

minimize $\|R(x'_1) - R(x_2)\|_2$

- Solve via gradient decent
 - here: projected gradient decent (PGD), similar to the attack
 - take a step toward the objective, project back on the region defined step size
 - no overall projection on a constraint (as in the attack)
- Expectation: x_1 ' is visually similar to x_2 (no constraint on δ)

Counterexample

find
$$x'_1 \coloneqq x_1 + \delta$$

minimize $\|R(x'_1) - R(x_2)\|_2$



But on a robustly trained network

find **minimize**

 $\begin{aligned} x_1' &\coloneqq x_1 + \delta \\ \|R(x_1') - R(x_2)\|_2 \end{aligned}$

Robust Net: trained with PGD.



Implications for Feature Visualization

- Feature Visualization by optimization requires heavy regularization
 - to encode "looks like an image" prior
 - else optimization produces noise

• What happens without regularization?

Implications for Feature Visualization

find $\boldsymbol{\chi}$ $score(x) - \lambda_1 R_1(x) - \dots - \lambda_m R_m(x)$ maximize $score(x) = mean(layer_n[x, y, z])$ where

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



layer_[x,y,z]

layer_[:,:,z]

 $layer_{1}[:,:,:]^{2}$



softmax[k]

Implications for Feature Visualization



Maximizing different coordinates (i)



Engstrom et al., 2019

Implications for Feature Attribution

- On non-robust networks gradient based-attribution methods require condition to show a "clean" saliency map
- On robust networks the gradient is better is aligned with human expectation (as the features are better aligned)

airplane

froa





(a) MNIST

(b) CIFAR-10

(c) Restricted ImageNet

Gradient of classfication loss w.r.t input

Why?

- Robust Neural Networks rely on different features
 - these features are aligned with human perception
 - e.g. these features are robust to small perturbations
 - thus visualizations are what we expect
- Non-Robust Neural Networks learn non-robust spurious features in the data
 - optimization shows these
 - but additional regularization hides them
 - these features may allow for higher accuracy

Consequence: Image Manipulation

The representation learned by the robust classifier is robust and versatile enough to allow for various computer vision tasks just by optimizing w.r.t. the input..



Santurkar et al., 2019

Why do adversarial examples exist?

 Recall (lecture 2): Neural Networks are too linear



• Additional Reason:

Non-Robust Neural Networks learn features that statistically are only weakly connected with the input

Lecture Summary

Feature Visualization & Feature Attribution

Connection to adversarial Robustness

Source (x_1)



Target (x_2)









Standard (x'_1) Robust (x'_1)