Lecture 8: Interpretability

Ranjay Krishna



Administrative: Assignment 2

Due 4/27 (extended) 11:59pm

- Multi-layer Neural Networks,
- Image Features,
- Optimizers

Ranjay Krishna

Administrative: Assignment 3

Due 5/11 11:59pm

- Normalization Layers,
- Dropout,
- CNNs

Will be released tomorrow

Ranjay Krishna



Administrative: Fridays

This Friday

Quantization





Administrative: Course Project

Project proposal due 4/29 11:59pm

Come to office hours to talk about your ideas





Visualizing and Understanding





Today: What's going on inside ConvNets?

This image is CC0 public domain



dense 192 128 2048 2048 192 48 128 224 dense densé TT 1000 192 192 128 Max 2048 2048 pooling Max Max 128 pooling pooling ŢŢŢŢ

Class Scores: 1000 numbers

Input Image: 3 x 224 x 224

What are the intermediate features looking for?

Krizhevsky et al, "ImageNet Classification with Deep Convolutional Neural Networks", NIPS 2012. Figure reproduced with permission.

Ranjay Krishna

Lecture 8 - 7

Today's agenda

Visualizing what models have learned:

- Visualizing filters
- Visualizing final layer features
- Visualizing activations

Understanding input pixels

- Identifying important pixels
- Saliency via backprop
- Guided backprop to generate images
- Gradient ascent to visualize features

Adversarial perturbations Concept Vectors

Ranjay Krishna

Lecture 8 - 8

Today's agenda

Visualizing what models have learned:

- Visualizing filters
- Visualizing final layer features
- Visualizing activations

Understanding input pixels

- Identifying important pixels
- Saliency via backprop
- Guided backprop to generate images
- Gradient ascent to visualize features

Adversarial perturbations Concept Vectors

Ranjay Krishna

Lecture 8 - 9

Interpreting a Linear Classifier: Visual Viewpoint







Ranjay Krishna

Lecture 8 - 10

First Layer: Visualize Filters



AlexNet: 64 x 3 x 11 x 11

Krizhevsky, "One weird trick for parallelizing convolutional neural networks", arXiv 2014 He et al, "Deep Residual Learning for Image Recognition", CVPR 2016 Huang et al, "Densely Connected Convolutional Networks", CVPR 2017

Ranjay Krishna

Lecture 8 - 11



First Layer: Visualize Filters



Max Max pooling Ma

AlexNet: 64 x 3 x 11 x 11

Krizhevsky, "One weird trick for parallelizing convolutional neural networks", arXiv 2014 He et al, "Deep Residual Learning for Image Recognition", CVPR 2016 Huang et al, "Densely Connected Convolutional Networks", CVPR 2017

Ranjay Krishna

Lecture 8 - 12

Visualize the filters/kernels (raw weights)

We can visualize filters at higher layers, but not that interesting Weights: 医骨骨 网络白色 化化化合金 化合金合金 Weights: (我们是我们的我们的我们的你的。(你们是我们的我们的你们的我们的我们的?" 如果我想出来的你们(我们会没有非常知道的自己的过去)(我们都能能能能能能能能能能。 · 模試)(並当由目書語並將此約定系書書方的)(物和指面與問時展開目標書種建成種)(書包書記書 而發行的內部保護保護部(保護局局等部所與保護部合管理的)(開設市務運動局部的運動 調査整備)(開始整約編集整整編装整備の開設)(は非計算は印刷書用計算用的計算)(構成構 那些這個這些可能是非常的。(你是這個的是是不是這個是是是有意思)(美力的水量是我们在非 要要到內面部(國際國法会局部建築局部的部分)(在空空的全面的和市市市市市市市)(国 医原油学会会 医原氨基氏试验 (口名名法律法法法口法法法法法法) (法以承知法法法法 新聞型化活動可**用**對) Weights: (这是当然能够建筑的萨思建筑基本的社会)(他是这些老师是是是是我的中国生活的。

)(國際軍業局部電気局通知型機械成長務委員會)(局理法院和部長局制務を受益局等なる)

20 x 16 x 7 x 7

layer 2 weights

layer 1 weights

16 x 3 x 7 x 7

Ranjay Krishna

Lecture 8 - 13

Last Layer

FC7 layer



4096-dimensional feature vector for an image (layer immediately before the classifier)

Run the network on many images, collect the feature vectors

Ranjay Krishna

Lecture 8 - 14

Last Layer: Nearest Neighbors

4096-dim vector

Test image L2 Nearest neighbors in feature space





Krizhevsky et al, "ImageNet Classification with Deep Convolutional Neural Networks", NIPS 2012. Figures reproduced with permission.

Ranjay Krishna

Lecture 8 - 15

April 24, 2025

Recall: Nearest neighbors in <u>pixel</u> space

Last Layer: Learned Metric for "Semantic" Search



Ranjay Krishna

Lecture 8 - 16

Last Layer: Modern Day Search



coactive.ai

Ranjay Krishna

Last Layer: Dimensionality Reduction

Visualize the "space" of FC7 feature vectors by reducing dimensionality of vectors from 4096 to 2 dimensions

Simple algorithm: Principal Component Analysis (PCA)

More complex: **t-SNE**





Van der Maaten and Hinton, "Visualizing Data using t-SNE", JMLR 2008 Figure copyright Laurens van der Maaten and Geoff Hinton, 2008. Reproduced with permission.

Ranjay Krishna

Lecture 8 - 18

Last Layer: Dimensionality Reduction



Van der Maaten and Hinton, "Visualizing Data using t-SNE", JMLR 2008 Krizhevsky et al, "ImageNet Classification with Deep Convolutional Neural Networks", NIPS 2012. Figure reproduced with permission.





See high-resolution versions at http://cs.stanford.edu/people/karpathy/cnnembed/

Ranjay Krishna

Lecture 8 - 19

Visualizing Activations

https://www.youtube.com/watch?v=AgkflQ4IGaM

conv5 feature map is 128x13x13; visualize as 128 13x13 grayscale images



Yosinski et al, "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014. Figure copyright Jason Yosinski, 2014. Reproduced with permission.

Ranjay Krishna

Lecture 8 - 20

Today's agenda

Visualizing what models have learned:

- Visualizing filters
- Visualizing final layer features
- Visualizing activations

Understanding input pixels

- Identifying important pixels
- Saliency via backprop
- Guided backprop to generate images
- Gradient ascent to visualize features

Adversarial perturbations Concept Vectors

Ranjay Krishna

Lecture 8 - 21

Maximally Activating Patches





Pick a layer and a channel; e.g. conv5 is 128 x 13 x 13, pick channel 17/128

Run many images through the network, record values of chosen channel

Visualize image patches that correspond to maximal activations



Springenberg et al, "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015 Figure copyright Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller, 2015; reproduced with permission.

Ranjay Krishna

Lecture 8 - 22

Which pixels matter: Saliency via Occlusion

Mask part of the image before feeding to CNN, check how much predicted probabilities change





P(elephant) = 0.95





P(elephant) = 0.75

Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014

Boat image is CC0 public domain Elephant image is CC0 public domain Go-Karts image is CC0 public domain

Ranjay Krishna

Lecture 8 - 23

Which pixels matter: Saliency via Occlusion

Mask part of the image before feeding to CNN, check how much predicted probabilities change









Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014

Boat image is CC0 public domain Elephant image is CC0 public domain Go-Karts image is CC0 public domain



African elephant, Loxodonta africana



April 24, 2025

- 0.8

Ranjay Krishna

Saliency via Occlusion: Shapley Values







$$P(corgi) = 0.8$$



Credit: Ian Covert; Lundberg & Lee 2017

April 24, 2025

Ranjay Krishna

Which pixels matter: Saliency via Backprop

Forward pass: Compute probabilities



Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

Ranjay Krishna

Lecture 8 - 26

Which pixels matter: Saliency via Backprop

Forward pass: Compute probabilities



Compute gradient of (unnormalized) class score with respect to image pixels, take absolute value and max over RGB channels

Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.



Ranjay Krishna

Lecture 8 - 27

Saliency Maps



Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

Ranjay Krishna

Lecture 8 - 28

Saliency Maps: Segmentation without supervision



Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission. Rother et al, "Grabcut: Interactive foreground extraction using iterated graph cuts", ACM TOG 2004

Ranjay Krishna

Use GrabCut on

saliency map

Lecture 8 - 29

Saliency maps: Uncovers biases

Such methods also find biases

wolf vs dog classifier looks is actually a snow vs no-snow classifier



(a) Husky classified as wolf



(b) Explanation

Figures copyright Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin, 2016; reproduced with permission. Ribeiro et al, ""Why Should I Trust You?" Explaining the Predictions of Any Classifier", ACM KDD 2016

Ranjay Krishna



Intermediate Features via (guided) backprop



Pick a single intermediate channel, e.g. one value in 128 x 13 x 13 conv5 feature map

Compute gradient of activation value with respect to image pixels

Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014 Springenberg et al, "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015

Ranjay Krishna

Lecture 8 - 31

Intermediate Features via (guided) backprop





b) Forward pass

1	-1	5		1	0	5
2	-5	-7	\rightarrow	2	0	0
-3	2	4		0	2	4
	_	_			_	_

Rel U

Backward pass: backpropagation

	0	2	4
		_	_
	-2	3	-1
←	6	-3	1
	2	-1	3

Pick a single intermediate neuron, e.g. one value in 128 x 13 x 13 conv5 feature map

Compute gradient of neuron value with respect to image pixels

Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014 Springenberg et al, "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015 Backward pass: guided backpropagation

0	0		-2	3	-1
0	0	->	6	-3	1
0	3		2	-1	3

Images come out nicer if you only backprop positive gradients through each ReLU (guided backprop)

Figure copyright Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller, 2015; reproduced with permission.

6

0

Ranjay Krishna

Lecture 8 - 32

Intermediate features via (guided) backprop





Guided Backprop

Maximally activating patches (Each row is a different neuron)

Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014 Springenberg et al, "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015 Figure copyright Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller, 2015; reproduced with permission.

Ranjay Krishna



Intermediate features via (guided) backprop



Guided Backprop

April 24, 2025

Maximally activating patches (Each row is a different neuron)

Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014 Springenberg et al, "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015 Figure copyright Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller, 2015; reproduced with permission.

Ranjay Krishna

Class Activation Mapping (CAM)



Ranjay Krishna

Lecture 8 - 35

Class Activation Mapping (CAM)



$$F_k = \frac{1}{HW} \sum_{hw} f_{h,w,k}$$

Ranjay Krishna

Lecture 8 - 36


Ranjay Krishna

Lecture 8 - 37



Ranjay Krishna

Lecture 8 - 38



Ranjay Krishna

Lecture 8 - 39



Ranjay Krishna

Lecture 8 - 40



Class activation maps of top 5 predictions



Class activation maps for one object class

Zhou et al, "Learning Deep Features for Discriminative Localization", CVPR 2016

Ranjay Krishna

Lecture 8 - 413

Problem: Can only apply to last conv



Class activation maps of top 5 predictions



Class activation maps for one object class

Ranjay Krishna

Lecture 8 - 424

1. Pick any layer, with activations $A \in \mathbb{R}^{H \times W \times K}$

Selvaraju et al, "Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization", CVPR 2017

Ranjay Krishna

Lecture 8 - 43⁵

- 1. Pick any layer, with activations $A \in \mathbb{R}^{H \times W \times K}$
- 2. Compute gradient of class score S_{C} with respect to A:

$$\frac{\partial S_c}{\partial A} \in \mathbb{R}^{H \times W \times K}$$

Selvaraju et al, "Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization", CVPR 2017

Ranjay Krishna



1. Pick any layer, with activations $A \in \mathbb{R}^{H \times W \times K}$

2. Compute gradient of class score S_c with respect to A:

 $\frac{\partial S_c}{\partial A} \in \mathbb{R}^{H \times W \times K}$

3. Global Average Pool the gradients to get weights $\alpha \in \mathbb{R}^{K}$:

$$\alpha_{k} = \frac{1}{HW} \sum_{h,w} \frac{\partial S_{c}}{\partial A_{h,w,k}}$$

Selvaraju et al, "Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization", CVPR 2017

Ranjay Krishna

Lecture 8 - 45⁵



1. Pick any layer, with activations $A \in \mathbb{R}^{H \times W \times K}$

2. Compute gradient of class score S_c with respect to A:

$$\frac{\partial S_c}{\partial A} \in \mathbb{R}^{H \times W \times K}$$

3. Global Average Pool the gradients to get weights $\alpha \in \mathbb{R}^{K}$:

$$\alpha_{k} = \frac{1}{HW} \sum_{h,w} \frac{\partial S_{c}}{\partial A_{h,w,k}}$$

4. Compute activation map $M^c \in \mathbb{R}^{H,W}$:

$$M_{h,w}^{c} = ReLU\left(\sum_{k} \alpha_{k} A_{h,w,k}\right)$$

Selvaraju et al, "Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization", CVPR 2017

Ranjay Krishna

Lecture 8 - 465



Selvaraju et al, "Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization", CVPR 2017

Ranjay Krishna

Lecture 8 - 479

Can also be applied beyond classification models, e.g. image captioning



A group of people flying kites on a beach

A man is sitting at a table with a pizza

Selvaraju et al, "Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization", CVPR 2017

Ranjay Krishna

Lecture 8 - 48°

(Guided) backprop:

Find the part of an image that a neuron responds to

Ranjay Krishna

Gradient ascent:

Generate a synthetic image that maximally activates a neuron

April 24, 2025



1. Initialize image to zeros

$$\arg\max_{I} S_c(I) - \lambda \|I\|_2^2$$

score for class c (before Softmax)

April 24, 2025

Repeat:

- 2. Forward image to compute current scores
- 3. Backprop to get gradient of neuron value with respect to image pixels
- 4. Make a small update to the image

Ranjay Krishna



$$\arg\max_{I} S_c(I) - \lambda \|I\|_2^2$$

Simple regularizer: Penalize L2 norm of generated image

Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014. Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

Ranjay Krishna

Lecture 8 - 51

-51



$$\arg\max_{I} S_c(I) - \lambda \|I\|_2^2$$

Simple regularizer: Penalize L2 norm of generated image



dumbbell

cup







bell pepper

lemon

husky

April 24, 2025

Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

Ranjay Krishna



$$\arg\max_{I} S_c(I) - \lambda \|I\|_2^2$$

Simple regularizer: Penalize L2 norm of generated image



Yosinski et al, "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014. Figure copyright Jason Yosinski, Jeff Clune, Anh Nguyen, Thomas Fuchs, and Hod Lipson, 2014. Reproduced with permission. goose

ostrich

limousine

April 24, 2025

Ranjay Krishna

Lecture 8 - 53

$$\arg\max_{I} S_c(I) - \lambda \|I\|_2^2$$

Better regularizer: Penalize L2 norm of image; also during optimization periodically

- (1) Gaussian blur image
- (2) Clip pixels with small values to 0
- (3) Clip pixels with small gradients to 0

Yosinski et al, "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014.

Ranjay Krishna

Lecture 8 - 54



$$\arg\max_{I} S_c(I) - \lambda \|I\|_2^2$$

Better regularizer: Penalize L2 norm of image; also during optimization periodically

- (1) Gaussian blur image
- (2) Clip pixels with small values to 0
- (3) Clip pixels with small gradients to 0



Flamingo



Ground Beetle



Pelican



Indian Cobra

April 24, 2025

Yosinski et al, "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014. Figure copyright Jason Yosinski, Jeff Clune, Anh Nguyen, Thomas Fuchs, and Hod Lipson, 2014. Reproduced with permission.

Ranjay Krishna

Lecture 8 - 55

$$\arg\max_{I} S_c(I) - \lambda \|I\|_2^2$$

Better regularizer: Penalize L2 norm of image; also during optimization periodically

- (1) Gaussian blur image
- (2) Clip pixels with small values to 0
- (3) Clip pixels with small gradients to 0



Hartebeest



Station Wagon



Billiard Table



Black Swan

April 24, 2025

Yosinski et al, "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014. Figure copyright Jason Yosinski, Jeff Clune, Anh Nguyen, Thomas Fuchs, and Hod Lipson, 2014. Reproduced with permission.

Ranjay Krishna

Lecture 8 - 56

Use the same approach to visualize intermediate features



Lecture 8 - 57

Yosinski et al, "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014. Figure copyright Jason Yosinski, Jeff Clune, Anh Nguyen, Thomas Fuchs, and Hod Lipson, 2014. Reproduced with permission.

Ranjay Krishna

57

Adding "multi-faceted" visualization gives even nicer results: (Plus more careful regularization, center-bias)

Reconstructions of multiple feature types (facets) recognized by the same "grocery store" neuron



Corresponding example training set images recognized by the same neuron as in the "grocery store" class



April 24, 2025

Nguyen et al, "Multifaceted Feature Visualization: Uncovering the Different Types of Features Learned By Each Neuron in Deep Neural Networks", ICML Visualization for Deep Learning Workshop 2016. Figures copyright Anh Nguyen, Jason Yosinski, and Jeff Clune, 2016; reproduced with permission.

Lecture 8 - 58

Ranjay Krishna



Lecture 8 - 59

Nguyen et al, "Multifaceted Feature Visualization: Uncovering the Different Types of Features Learned By Each Neuron in Deep Neural Networks", ICML Visualization for Deep Learning Workshop 2016. Figures copyright Anh Nguyen, Jason Yosinski, and Jeff Clune, 2016; reproduced with permission.

Ranjay Krishna

59

Optimize in FC6 latent space instead of pixel space:



Nguyen et al, "Synthesizing the preferred inputs for neurons in neural networks via deep generator networks," NIPS 2016 Figure copyright Nguyen et al, 2016; reproduced with permission.

Ranjay Krishna

Lecture 8 - 60

60

Optimize in FC6 latent space instead of pixel space:



Nguyen et al, "Synthesizing the preferred inputs for neurons in neural networks via deep generator networks," NIPS 2016 Figure copyright Nguyen et al, 2016; reproduced with permission.

Ranjay Krishna

Lecture 8 - 61

61

Today's agenda

Visualizing what models have learned:

- Visualizing filters
- Visualizing final layer features
- Visualizing activations

Understanding input pixels

- Identifying important pixels
- Saliency via backprop
- Guided backprop to generate images
- Gradient ascent to visualize features

Adversarial perturbations

Concept Vectors

Ranjay Krishna

Lecture 8 - 62

Fooling Images / Adversarial Examples

- (1) Start from an arbitrary image
- (2) Pick an arbitrary incorrect class
- (3) Modify the image to maximize the class
- (4) Repeat until network is fooled

Fooling Images / Adversarial Examples

African elephant





iPod



Difference



schooner









10x Difference



Boat image is CC0 public domain Elephant image is CC0 public domain

Ranjay Krishna

Lecture 8 - 64



Fooling Images / Adversarial Examples

African elephant





iPod



Difference



schooner







10x Difference



April 24, 2025

Check out lan Goodfellow's lecture from 2017

Boat image is CC0 public domain Elephant image is CC0 public domain

Ranjay Krishna

Lecture 8 - 65

Fooling Person Detectors and Self-driving Cars









frame 120

frame 150









Lecture 8 - 66

Xu et al., 2019; Eykholt et al., 2018



Fooling Images / Adversarial Exa

Universal perturbations

Moosavi-Dezfooli, Seyed-Mohsen, et al. "Universal adversarial perturbations." Proceedings of the IEEE conference on computer vision and pattern recognition. 2017.

Figure reproduced with permission



Ranjay Krishna

Feature Inversion

Given a CNN feature vector for an image, find a new image that:

- Matches the given feature vector -
- "looks natural" (image prior regularization) -

$$\mathbf{x}^{*} = \underset{\mathbf{x} \in \mathbb{R}^{H \times W \times C}}{\operatorname{argmin}} \ell(\Phi(\mathbf{x}), \Phi_{0}) + \lambda \mathcal{R}(\mathbf{x})} \qquad \text{vector}$$

$$\ell(\Phi(\mathbf{x}), \Phi_{0}) = \|\Phi(\mathbf{x}) - \Phi_{0}\|^{2} \qquad \text{image}$$

$$\mathcal{R}_{V^{\beta}}(\mathbf{x}) = \sum_{i,j} \left((x_{i,j+1} - x_{ij})^{2} + (x_{i+1,j} - x_{ij})^{2} \right)^{\frac{\beta}{2}} \qquad \text{Total Variation regularizer}$$

$$(\text{encourages spatial}$$

$$(\text{an and Vedaldi, "Understanding Deep Image Representations by Inverting Them", CVPR 2015} \qquad \text{Solven feature}$$

April 24, 2025

Mahendran and Vedaldi, "Understanding Deep Image Representations by Inverting Them", CVPR 2015

Ranjay Krishna

Feature Inversion

Reconstructing from different layers of VGG-16



Mahendran and Vedaldi, "Understanding Deep Image Representations by Inverting Them", CVPR

Ranjay Krishna

Lecture 8 - 69

DeepDream: Amplify Existing Features

Rather than synthesizing an image to maximize a specific neuron, instead try to **amplify** the neuron activations at some layer in the network





Choose an image and a layer in a CNN; repeat:

- 1. Forward: compute activations at chosen layer
- 2. Set gradient of chosen layer equal to its activation
- 3. Backward: Compute gradient on image
- 4. Update image

Mordvintsev, Olah, and Tyka, "Inceptionism: Going Deeper into Neural Networks", <u>Google Research Blog</u>. Images are licensed under <u>CC-BY 4.0</u>

ril 24, 2025

Ranjay Krishna

DeepDream: Amplify Existing Features

Rather than synthesizing an image to maximize a specific neuron, instead try to **amplify** the neuron activations at some layer in the network





Choose an image and a layer in a CNN; repeat:

- 1. Forward: compute activations at chosen layer
- 2. Set gradient of chosen layer *equal to its activation*
- 3. Backward: Compute gradient on image
- 4. Update image

Ranjay Krishna

Lecture 8 - 71

Equivalent to: $I^* = \arg \max_{I} \sum_{i} f_i(I)^2$

Mordvintsev, Olah, and Tyka, "Inceptionism: Going Deeper into Neural Networks", <u>Google Research Blog</u>. Images are licensed under <u>CC-BY 4.0</u>

il 24, 2025



Ranjay Krishna






Lecture 8 - 74



Lecture 8 - 75







Image is licensed under <u>CC-BY 4.0</u>

April 24, 2025

Ranjay Krishna

Today's agenda

Visualizing what models have learned:

- Visualizing filters
- Visualizing final layer features
- Visualizing activations

Understanding input pixels

- Identifying important pixels
- Saliency via backprop
- Guided backprop to generate images
- Gradient ascent to visualize features

Adversarial perturbations

Concept Vectors

Ranjay Krishna

Lecture 8 - 78

Concept activation vectors



Let's see if a neural network has learned a specific concept and uses it effectively.

Example use case: **Q1.** Has it learnt what stripes are? **Q2.** Can it identify the category "zebra" by using the concept "stripes"

April 24, 2025

Kim, Been, et al. "Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav)." ICML, 2018.

Ranjay Krishna

Concept activation vectors



Kim, Been, et al. "Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav)." ICML, 2018.

Ranjay Krishna

Lecture 8 - 80

Calculate if the gradient for that layer when predicting "zebra" matches the classifier



Kim, Been, et al. "Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav)." ICML, 2018.

Ranjay Krishna

Lecture 8 - 81

April 24, 2<u>025</u>



Many methods for understanding CNN representations

Activations: Nearest neighbors, Dimensionality reduction, maximal patches, occlusion Gradients: Saliency maps, class visualization, fooling images, feature inversion

Adversarial Examples: To confuse the models Concept Vectors: Human interpretable probing method

Ran	jay	Kris	hna

Lecture 8 - 82

82

Next time:

Introduction to Language



Lecture 8 - 83

83

