Generating Images from CNNs

Designing, Visualizing and Understanding Deep Neural Networks

CS W182/282A

Instructor: Sergey Levine UC Berkeley



What does the brain see?









What do convolutional networks "see"?



- > Interpret what neural nets pay attention to
- Understand when they'll work and when they won't
- Compare different models and architectures
- Manipulate images based on conv net responses





What do we visualize?

Option 1: visualize the **filter** itself



Option 2: visualize stimuli that activate a "neuron"



Visualizing filters





Can't really visualize higher layers in a way that makes much sense

Visualizing neuron responses

Idea 1: look for images that maximally "excite" specific units



Visualizing neuron responses



Figure 4: Top regions for six pool₅ units. Receptive fields and activation values are drawn in white. Some units are aligned to concepts, such as people (row 1) or text (4). Other units capture texture and material properties, such as dot arrays (2) and specular reflections (6).

Girshick, Donahue, Darrell, Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. 2014

Idea 1: See which image pixels maximally influence the value at some unit

what does "influence" mean?

given a pixel x_{ij} and a unit a_{mnp}^{ℓ} how much does changing x_{ij} change a_{mnp}^{ℓ} ?

 $\frac{da_{mnp}^{\ell}}{dx_{ij}} \quad \text{how do we get this quantity?} \\ \text{backpropagation!}$

1. set δ to be same size as a^{ℓ} 2. set $\delta_{mnp} = 1$, all other entries to 0 3. backprop from layer ℓ to the image 4. last δ gives us $\frac{da_{mnp}^{\ell}}{dx}$



forward pass: calculate each $a^{(i)}$ and $z^{(i)}$ backward pass: $\frac{da_{mnp}^{\ell}}{da^{\ell}} \leftarrow$ zero in each position except mnp (which is 1) initialize $\delta = \frac{d\mathcal{L}}{dz} + \frac{da_{mnp}^{\ell}}{da^{\ell}} \leftarrow$ zero in each position except mnp (which is 1) for each f with input x_f & params θ_f from end to start: $\frac{d\mathcal{L}}{d\theta_f} \leftarrow \frac{df}{d\theta_f} \delta$ $\delta \leftarrow \frac{df}{dx_f} \delta$

Idea 1: See which image pixels maximally influence the value at some unit



basic idea behind "guided backpropagation":

 Just using backprop is not very interpretable, because many units in the network contribute positive or negative gradients

 da_{mnp}^{ℓ}

 dx_{ij}

If we keep just the **positive** gradients, we'll avoid some of the complicated negative contributions, and get a cleaner signal

heuristically change backward step in ReLU:

for each entry δ_{ijk} , set it to max $(0, \delta_{ijk})$ "zero out negative gradients at each layer"



slightly modified backprop gives us this:

"guided backpropagation"



Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller. Striving for Simplicity: The All Convolutional Net. 2014.

Idea 1: See which image pixels maximally influence the value at some unit



Figure 1: Schematic of visualizing the activations of high layer neurons. a) Given an input image, we perform the forward pass to the layer we are interested in, then set to zero all activations except one and propagate back to the image to get a reconstruction. b) Different methods of propagating back through a ReLU nonlinearity. c) Formal definition of different methods for propagating a output activation *out* back through a ReLU unit in layer *l*; note that the 'deconvnet' approach and guided backpropagation do not compute a true gradient but rather an imputed version.



basic idea behind "guided backpropagation":

- Just using backprop is not very interpretable, because many units in the network contribute positive or negative gradients
- If we keep just the **positive** gradients, we'll avoid some of the complicated negative contributions, and get a cleaner signal

heuristically change backward step in ReLU:

for each entry δ_{ijk} , set it to max $(0, \delta_{ijk})$ "zero out negative gradients at each layer"

Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller. Striving for Simplicity: The All Convolutional Net. 2014.

Idea 1: See which image pixels maximally influence the value at some unit



Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller. Striving for Simplicity: The All Convolutional Net. 2014.

Visualizing Features with Backpropagation

Idea 2: "Optimize" the image to maximally activate a particular unit

 $\underline{da_{mnp}^{\ell}}$

Before: just compute



what optimization problem is this solving?



Idea 2: "Optimize" the image to maximally activate a particular unit



Simonyan, Vidaldi, Zisserman. Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps. 2013.

Visualizing "classes"

 $x \leftarrow \arg \max_{x} S(x) + R(x)$ \uparrow simple choice: $R(x) = \lambda ||x||_{2}^{2}$

Which unit to maximize?

Let's try maximizing class labels first

Important detail: not maximizing class probabilities, but the activations right before the softmax





Simonyan, Vidaldi, Zisserman. Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps. 2013.

Visualizing "classes"

 $x \leftarrow \arg\max_{x} S(x) + R(x)$ fslightly more nuanced regularizer

- 1. Update image with gradient
- 2. Blur the image a little
- 3. Zero out any pixel with small value
- 4. Repeat



Flamingo

Ground Beetle



Pelican



Indian Cobra



Hartebeest



Station Wagon



Billiard Table



Black Swan

Slide based on CS231n, Fei-Fei Li & Andrej Karpathy

Visualizing units

 $x \leftarrow \arg\max_x S(x) + R(x)$ fslightly more nuanced regularizer

- 1. Update image with gradient
- 2. Blur the image a little
- 3. Zero out any pixel with small value
- 4. Repeat



Slide based on CS231n, Fei-Fei Li & Andrej Karpathy

Visualizing units



- 1. Update image with gradient
- 2. Blur the image a little
- 3. Zero out any pixel with small value
- 4. Repeat

Slide based on CS231n, Fei-Fei Li & Andrej Karpathy

Visualizing units



- 1. Update image with gradient
- 2. Blur the image a little
- 3. Zero out any pixel with small value
- 4. Repeat



Slide based on CS231n, Fei-Fei Li & Andrej Karpathy

Deep Dream & Style Transfer

Using backprop to modify pictures?

Before:



Now:



What is channel 17 in layer conv5 looking at?

What would it look like if it were a Monet painting?



How are these questions related?

Intuition: Monet paintings have particular feature distributions, we can "transport" these distributions to other images! ??

Looking for patterns in clouds







making the eyes more blue will make this more cat-like





inception_4c/output



1. Pick a layer

- 2. Run forward pass to compute activations at that layer
- 3. Set delta to be **equal** to the activations
- 4. Backprop and apply the gradient
- 5. Repeat



inception_4c/output







Remember this:

$$x \leftarrow \arg\max_x S(x) + R(x)$$



DeepDream https://github.com/google/deepdream

Another idea

Another idea: instead of exaggerating the features in a single image, what if we make feature of one image look more like the features in **another**?



How do we quantify style?



Images from CS231n (Fei-Fei Li, Andrej Karpathy)

How do we quantify this?

How do we quantify style?



style

relationships between

features

e.g., lots of angles and straight lines

estimate by averaging over all positions in the image



```
feature covariance: \operatorname{Cov}_{km} = E[f_k f_m]
```

form Gram matrix: $G_{km} = \operatorname{Cov}_{km}$

If features have this covariance, then we have the right style

Surprising but true!

"each time we see a straight line, we also see this texture"

"each time we see a curve, we also see flat shading"

Images from CS231n (Fei-Fei Li, Andrej Karpathy) Gatys et al. A Neural Algorithm of Artist Style, 2015.

How do we quantify style?



style

relationships between

features

e.g., lots of angles and straight lines

feature covariance: $\operatorname{Cov}_{km} = E[f_k f_m]$ form Gram matrix: $G_{km} = \operatorname{Cov}_{km}$ \longleftarrow this comes from the style source image new image: $x \leftarrow \arg\min_x \mathcal{L}_{\operatorname{style}}(x) + \mathcal{L}_{\operatorname{content}}(x)$ G^{ℓ} : source image Gram matrix at layer ℓ $A^{\ell}(x)$: new image Gram matrix at layer ℓ $\mathcal{L}_{\operatorname{style}}(x) = \sum_{\ell} \sum_{km} \left(G_{km}^{\ell} - A_{km}^{\ell}(x) \right)^2 w_{\ell}$ Different weight on each layer, to prioritize relative contribution to (desired) style of different levels of abstraction

Images from CS231n (Fei-Fei Li, Andrej Karpathy) Gatys et al. A Neural Algorithm of Artist Style, 2015.

How do we quantify content?



5 features $rac{f_i}{}$

content

spatial positions of

features

e.g., where curves, edges, and corners are located

Pick specific layer for matching content

Just directly match the features!

$$\mathcal{L}_{\text{content}}(x) = \sum_{ij} \sum_{k} \left(f_{ijk}^{\ell}(x_{\text{content}}) - f_{ijk}^{\ell}(x) \right)^2$$

$$\mathcal{L}_{\text{style}}(x) = \sum_{\ell} \sum_{km} \left(G_{km}^{\ell} - A_{km}^{\ell}(x) \right)^2 w_{\ell}$$

new image: $x \leftarrow \arg \min_x \mathcal{L}_{style}(x) + \mathcal{L}_{content}(x)$

Images from CS231n (Fei-Fei Li, Andrej Karpathy) Gatys et al. A Neural Algorithm of Artist Style, 2015.

Style transfer

new image: $x \leftarrow \arg \min_x \mathcal{L}_{style}(x) + \mathcal{L}_{content}(x)$

 $\mathcal{L}_{\text{content}}(x) = \sum_{ij} \sum_{k} \left(f_{ijk}^{\ell}(x_{\text{content}}) - f_{ijk}^{\ell}(x) \right)^2$

$$\mathcal{L}_{\text{style}}(x) = \sum_{\ell} \sum_{km} \left(G_{km}^{\ell} - A_{km}^{\ell}(x) \right)^2 w_{\ell}$$



make your own easily on deepart.io





Gatys et al. A Neural Algorithm of Artist Style, 2015.