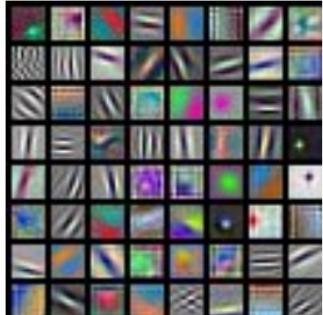


First layer of ConvNet:

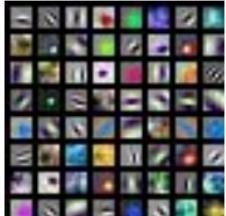
↳ Visualizing the learned filters/weights



AlexNet:
 $64 \times 3 \times 11 \times 11$



ResNet-18:
 $64 \times 3 \times 7 \times 7$



ResNet-101:
 $64 \times 3 \times 7 \times 7$



DenseNet-121:
 $64 \times 3 \times 7 \times 7$

↳ Rescaled to 0-255

Q. Why use filters for visualizing?

↳ Because the thing that causes the inner product to excite maximally is a copy of the thing you are taking an inner product with.

Second layers:

- ↳ Second layer receives these 16 7×7 input λ does 7×7 conv with 20 filters
- ↳ We can't intercept these because these are NOT connected to the input image but rather output of first layer.
- ↳ Also, SL learns the relationship in first layer that make more complex features

Last layers:

- ↳ If we use nearest neighbours but instead using pixels from the last layer
- ↳ The feature map of two differently oriented images end up close to each other.

t-distributed stochastic neighbour embeddings:

- ↳ Non linear dimension reduction method
- ↳ Use 4096 from layer λ apply t-SNE

- ↳ We have pixels of the image
- ↳ We have the 4096 d vectors

- ↳ We use t-SNE to get a 2 dimensional coordinate from 4096 vectors

NOTE:

1. A neuron corresponds to a single value that corresponds to a patch of the image
2. We use multiple filters to look at a multiple features

↳ We place the image at the coordinate

Maximally Activating patches:

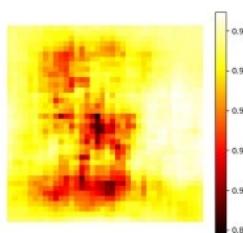
- ↳ Pick a layer ↳ a channel. Eg: conv5 is $128 \times 13 \times 13$, pick a channel
- ↳ Choose the patch of image that maximally activate it.
- ↳ This gives us a sense of what type of features that neuron is looking for

Occlusion experiment:

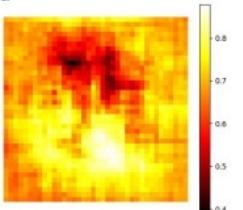
- ↳ Which parts of the input image leads the network classify it

1. Take an image
2. Block a region of the image
3. Replace it with mean pixel image
4. Go through the network
5. Go to step 2

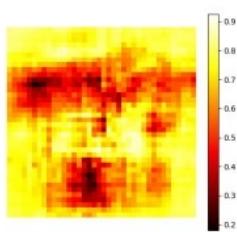
- ↳ If a blocked patch of the image caused the image to change drastically then that patch was really imp for the classification



African elephant, Loxodonta africana



high prob score corresponding to a blocked patch

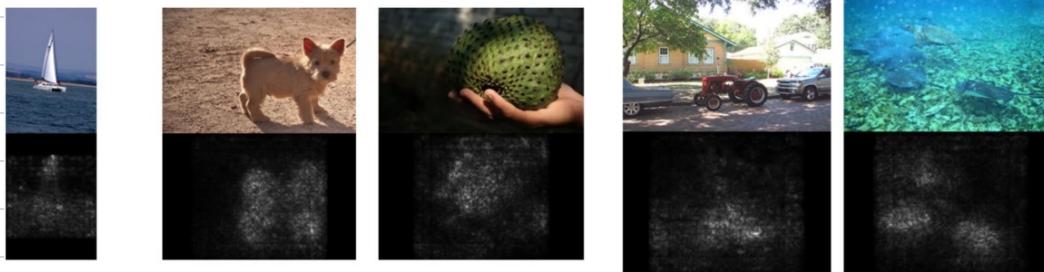


low prob score corresponding to a blocked patch

Saliency Maps:

↳ How to tell which pixel matters for classification?

↳ Take grad of predicted class score w.r.t the pixels of input image



Saliency maps for semantic segmentation:

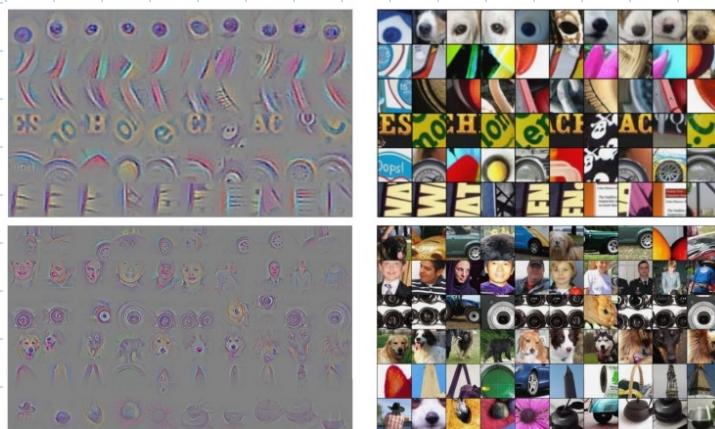
↳ Gradient segmentation algorithm + Saliency map ends up segmenting the image

↳ NOT very reliable

Intermediate features via (guided) backprop

↳ Take grad of nueron w.r.t the pixels of input image

↳ Images come out nice if you only backprop only through ten grads through each ReLU



Gradient Ascent:

↳ Weights remain same

↳ Change pixels to maximize the neuron activations

$$I^* = \operatorname{argmax}_I f(I) + R(I)$$

↓
Neuron
value

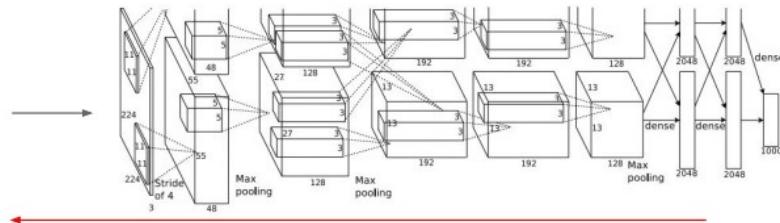
↳ Natural image
regularizer

$\|I\|_2^2$ norm of gen
image

$$\arg \max_I S_c(I) - \lambda \|I\|_2^2$$

score for class c (before Softmax)

1. Initialize image to zeros



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

↳ Instead of L_2 norm, we add

- ↳ gaussian blur image
- ↳ Clip pixels with small values to 0
- ↳ Clip pixels with small gradients to 0

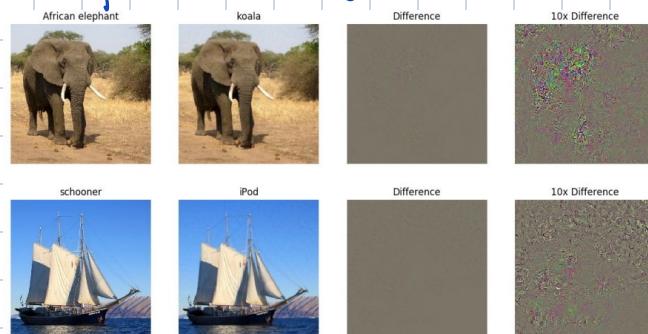
Another way:

↳ Maximize FC6 representation of an image rather than score

Fooling Images/Adversarial Examples

- ↳ Start from an arbitrary image → Elephant image
- ↳ Pick an arbitrary image → Koala bear
- ↳ Modify the image to maximize the class
- ↳ Repeat until the image is fooled

↳ Elephant would morph into koala bear but it does not?



DeepDream:

Choose an image & a layer in CNN:

- ↳ Forward: Compute activations at chosen layer
 - ↳ Set gradient of chosen layer equal to its activation
 - ↳ Backward: Compute gradient on image
 - ↳ Update image
- ↳ Amplify existing features

Feature Inversion:

- ↳ Run image through the net
- ↳ Record feature vector of that image
- ↳ Recreate image using that feature representation

$$\mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^{H \times W \times C}}{\operatorname{argmin}} \ell(\Phi(\mathbf{x}), \Phi_0) + \lambda \mathcal{R}(\mathbf{x})$$

Given feature vector

Features of new image

$$\ell(\Phi(\mathbf{x}), \Phi_0) = \|\Phi(\mathbf{x}) - \Phi_0\|^2 \rightarrow \text{minimize the dist b/w gen features & saved feature}$$
$$\mathcal{R}_{V^\beta}(\mathbf{x}) = \sum_{i,j} \left((x_{i,j+1} - x_{ij})^2 + (x_{i+1,j} - x_{ij})^2 \right)^{\frac{\beta}{2}}$$

Total Variation regularizer
(encourages spatial smoothness)

Texture synthesis:

- ↳ Given a small patch of some texture, generate larger image of same texture.

Regenerate pixels one at a time in scanline order; form neighbourhood of already generated pixels and copy nearest neighbour from input

- ↳ Works well for simple patterns

Neural texture Synthesis:

- ↳ Pass image through the network
- ↳ Take the activation volume
- ↳ Take two C -dimensional vectors
- ↳ Take outer product will give $C \times C$

↳ Do for all

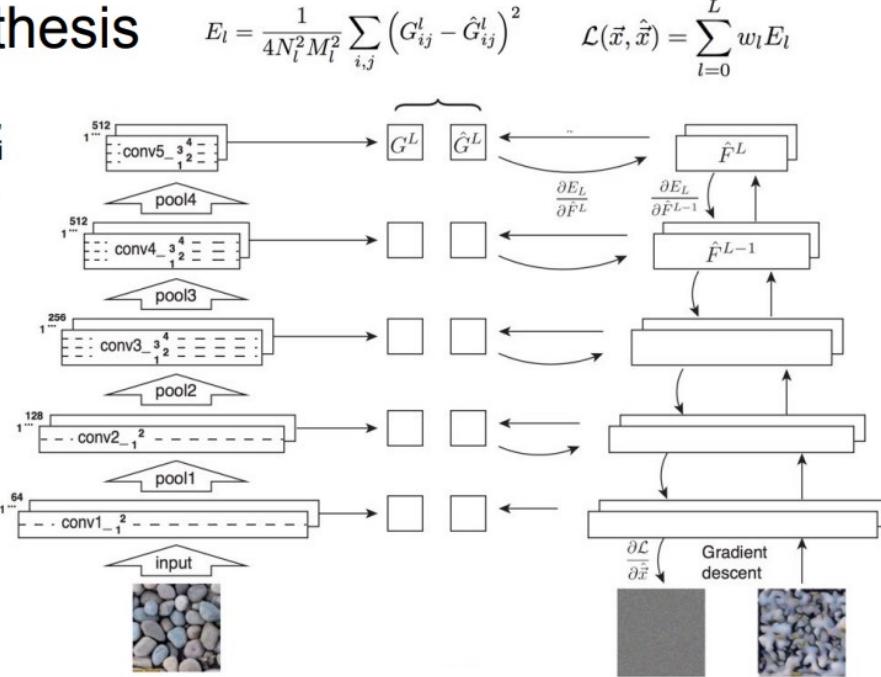
↳ Avg over $H \times W$ to get gram matrix

↳ Captures co-occurrence - i.e. if $i \neq j$ is large, the G_{ij} will be larger

↳ Recreate the gram matrix

Neural Texture Synthesis

1. Pretrain a CNN on ImageNet (VGG-19)
 2. Run input texture forward through CNN, record activations on every layer; layer i gives feature map of shape $C_i \times H_i \times W_i$
 3. At each layer compute the Gram matrix giving outer product of features:
- $$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l \text{ (shape } C_i \times C_i\text{)}$$
4. Initialize generated image from random noise
 5. Pass generated image through CNN, compute Gram matrix on each layer
 6. Compute loss: weighted sum of L2 distance between Gram matrices
 7. Backprop to get gradient on image
 8. Make gradient step on image
 9. GOTO 5



Style transfer:

Content image + Style image = style transfer

↳ What output will look like

↳ Texture

↳ Minimize feature represent & gram matrix loss

