Neural Style Transfer: Feature extractor, Training Loop

paper (https://arxiv.org/pdf/1508.06576.pdf)

The objective of Neural Style Transfer:

- ullet Given Content Image C
- ullet Given Style Image S
- ullet Create Generated Image G that is the Content image re-drawn in the "style" of the Style image



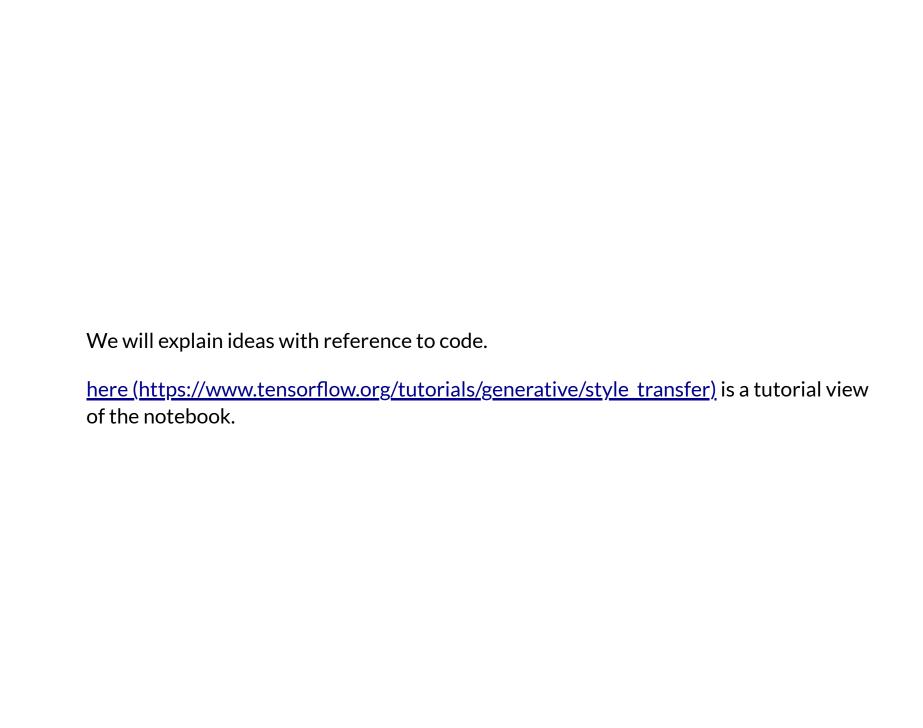




Style image S Content image C Generated image G

Neural Style Transfer highlights several themes we will encounter in the course

- The essential element of Deep Learning is
 - defining a Loss function that captures the semantics of the task
 - architecture is less important: just a tool
- The intermediate representations of a Deep Network has meaning, that can be leveraged
- Re-using existing models in novel ways leads to powerful results



Loss function: sum of Content Loss and Style Loss

We create a Loss Function $\mathcal{L}(C,G,S)$ and solve for the optimal generated image G^* $G^*=rgmin_G\mathcal{L}(C,G,S)$

where $\mathcal{L}(C,G,S)$ is the sum of

- ullet "content loss" $\mathcal{L}_{ ext{content}}$: dissimilarity of "content" of G and C
- ullet "style loss" $\mathcal{L}_{ ext{style}}$ dissimilarity of "style" of G and S

We solve for G^* using $rac{\partial \mathcal{L}}{\partial G}$

- ullet depending on how we write ${\cal L}$
 - minimize dissimilarity: Gradient Descent
 - maximize similarity; Gradient Ascent

That is: the "weights" we are optimizing are the pixels of image I.

How do we measure the dissimilarity of the "content"?

We can't just use plain MSE of the pixel-wise differences

• G is different than C, by definition (the "styles" are different)

And how do we define what the "style" of an image is?

• And how do we measure dissimilarity of the "style"?

We will use an alternate representation of each image

• such that we can compare alternate representations in a useful sense

The goal of using an alternate representation of an image

- is to capture the "semantics" (deeper, non-surface meaning) of an image
- rather than "syntax" (superficial surfrace meaning, literal pixels) of an image

_

Representation of intermediate layers as a measure of Style and Content

Recall that each layer in a multi-layer Neural Network is creating an alternate representation of the input.

Rather than directly comparing G with C (and G with S) our dissimilarity will be measured

- Not on raw images as seen by the human eye
- But on their alternate representations as created at some layer of a multi-layer
 Neural Network

That is: we will re-use a model $\mathbb C$ (e.g., VGG19)

- originally designed for Image Classification
- as a means to create two alternate representations of an image
- ullet one alternate representation of image I will encode the "content" of I
- ullet the other alternate representation of image I will encode the "style" of I

Suppose $\mathbb C$ consists of a sequence of CNN Layers

Let $\mathbb{C}_{(l)}$ denote the set of $n_{(l)}$ feature maps produced at layer l

- Feature map: value of one feature, at each spatial location
- $\mathbb{C}_{(l),j}$: feature map j

We choose

- One layer l_c of $\mathbb C$ and call it the "content representation" layer
 - Will tend to be shallow: closer to the input
 - Features of shallow layers will be more "syntax" than "semantics"
- ullet One layer l_s of ${\mathbb C}$ and call it the "style representation" layer
 - Will tend to be deep: closer to the output
 - Features of deep layers will be more "semantics" than "syntax"

For arbitrary image I, let

- ullet $\mathbb{C}_{(l_c)}(I)$
 - \blacksquare denote the feature maps of the Classifier $\mathbb C,$ on image I, at the "content representation" layer
- ullet $\mathbb{C}_{(l_s)}(I)$
 - \blacksquare denote the feature maps of the Classifier $\mathbb C$, on image I , at the "style representation" layer

Using the alternate representations derived from $\mathbb C$ we can define

- ullet $\mathcal{L}_{\mathrm{content}}$ as the dissimilarity of $\mathbb{C}_{(l_c)}(C)$ and $\mathbb{C}_{(l_c)}(G)$
- ullet $\mathcal{L}_{ ext{style}}$ as the dissimilarity of $\mathbb{C}_{(l_s)}(S)$ and $\mathbb{C}_{(l_c)}(G)$

Content Loss $\mathcal{L}_{content}$

We can now define the similarity of the "content" of Content Image ${\cal C}$ and "content" of Generated Image ${\cal G}$

- ullet by comparing $\mathbb{C}_{(l_c)}(C)$ and $\mathbb{C}_{(l_c)}(G)$
- ullet using sum of pixel-wise squared difference (MSE) of $\mathbb{C}_{(l_c)}(C)$ and $\mathbb{C}_{(l_c)}(G)$

Here is the code for content loss $\mathcal{L}_{\text{content}}$

- ullet base is $\mathbb{C}_{(l_c)}(C)$: the alternate representation of content image C
- ullet combination is $\mathbb{C}_{(l_c)}(G)$ the alternate representation of generated image G

```
def content_loss(base, combination):
    return tf.reduce_sum(tf.square(combination - base))
```

Style Loss $\mathcal{L}_{\mathrm{style}}$

Similarly, we can define the similarity of the "style" of Content Image ${\cal C}$ and "style" of Generated Image ${\cal G}$

ullet by comparing $\mathbb{C}_{(l_s)}(S)$ and $\mathbb{C}_{(l_s)}(G)$

For any image I: $\mathbb{C}_{(l)}(I)$ consists of $n_{(l)}$ feature maps.

We need to define what it means to compare $\mathbb{C}_{(l)}(I)$ and $\mathbb{C}_{(l)}(I')$.

The Gramm Matrix \mathbb{G} of $\mathbb{C}_{(l)}(I)$

- Has shape ($n_{(l)} imes n_{(l)}$)
- $ullet \ \mathbb{G}_{j,j'}(I) = \operatorname{correlation}(\operatorname{flatten}(\mathbb{C}_{(l),j}(I)), \ \operatorname{flatten}(\mathbb{C}_{(l),j'}(I)))$
 - lacktriangledown the correlation of the feature map j of $\mathbb{C}_{(l)}(I)$ with feature map j' of $\mathbb{C}_{(l)}(I')$

Intuitively, the Gramm Matrix

ullet measures the correlation of the values across pixel locations (flattened feature maps) of two feature maps of image I

Here is the code computing the "style loss" $\mathcal{L}_{ ext{style}}$

- ullet style is $\mathbb{C}_{(l_s)}(S)$: the alternate representation of style image S
- ullet combination is $\mathbb{C}_{(l_s)}(G)$ the alternate representation of generated image G

```
def style_loss(style, combination):
   S = gram_matrix(style)
   C = gram_matrix(combination)
   channels = 3
   size = img_nrows * img_ncols
   return tf.reduce_sum(tf.square(S - C)) / (4.0 * (channels**2) * (size**2))
```

Gradient ascent: generating $oldsymbol{G}$

We can find image G via Gradient Ascent

- Ascent versus Descent: we have measured *similarity* (correlation) rather than *dissimilarity*
- ullet Initialize G to noise
- Update pixel $G_{i,i',k}$ by $rac{\partial \mathcal{L}}{G_{i,i',k}}$

Feature extractor

One key coding trick that we will illustrate

ullet Obtaining the feature maps of the Classifier $\mathbb C$, on image I, at an arbitrary layer

We will call this tool the feature extractor

```
# Build a VGG19 model loaded with pre-trained ImageNet weights
model = vgg19.VGG19(weights="imagenet", include_top=False)

# Get the symbolic outputs of each "key" layer (we gave them unique names).
outputs_dict = dict([(layer.name, layer.output) for layer in model.layers])

# Set up a model that returns the activation values for every layer in
# VGG19 (as a dict).
feature_extractor = keras.Model(inputs=model.inputs, outputs=outputs_dict)
```

The feature_extractor code returns a dictionary

• mapping layer name to alternate representation at that layer

It is used within the compute_loss function

- to simultaneously compute (through threading on the first dimension) the alternate representations
- ullet of the three images C,S and G
- ullet base_image is the content image C
- ullet style_reference_image is the style image S
- ullet combination_image is the generated image G

```
def compute_loss(combination_image, base_image, style_reference_image):
    input_tensor = tf.concat(
        [base_image, style_reference_image, combination_image], axis=0
    )
    features = feature_extractor(input_tensor)
```

The alternate representations (at all layers) of each image

- is extracted from features later in the code
- here is the code using the alternate representations to eventually compute style loss
 - n.b., $\mathcal{L}_{\mathrm{sty[e}}$ is computed here over *several layers*
 - here is the code for one layer named layer_name

```
layer_features = features[layer_name]
style_reference_features = layer_features[1, :, :, :]
combination_features = layer_features[2, :, :, :]
```

Computing gradients

Here is the code to enable gradients to be computed

- tf.GradientTape records the forward pass
- to facilitate computation of gradients in the backward pass

You may want to review (from the Intro course)

- Back propagation: forward and backward pass (Training Neural Network Backprop.ipynb)
- <u>Computing analytic gradients in Keras</u>
 <u>(Training Neural Network Operation Forward and Backward Pass.ipynb)</u>

```
@tf.function
def compute_loss_and_grads(combination_image, base_image, style_reference_imag
e):
    with tf.GradientTape() as tape:
        loss = compute_loss(combination_image, base_image, style_reference_imag
e)
    grads = tape.gradient(loss, combination_image)
    return loss, grads
```

Recall the variables

- ullet base_image is the content image C
- ullet style_reference_image is the style image S
- ullet combination_image is the generated image G

Training loop

Here is the code for the training loop

- optimizer.apply_gradients([(grads, combination_image)])
- updates the "weights" combination_image using gradients grads

...

```
optimizer = keras.optimizers.SGD(
    keras.optimizers.schedules.ExponentialDecay(
        initial_learning_rate=100.0, decay_steps=100, decay_rate=0.96
    )
)
...
iterations = 4000
for i in range(1, iterations + 1):
    loss, grads = compute_loss_and_grads(
        combination_image, base_image, style_reference_image
    )
    optimizer.apply_gradients([(grads, combination_image)])
```

Notebooks links

<u>Full notebook (https://www.tensorflow.org/tutorials/generative/style_transfer)</u> is a tutorial view of the notebook we used for the code snippets.

```
In [2]: print("Done")
```

Done