

# CS231n Assignment 3

- RNN & LSTM image captioning
  - Network Visualization
    - o Saliency maps
    - o Class Visualization
    - o Fooling Images
  - Style Transfer
- Aside: Matrix Norm and Vector Norm
- Generative Adversarial Networks



# CS231n Assignment 3

function sample

feature ( $N, D$ )

captions ( $N, \text{max\_length}$ )

$\text{max\_length}$   
 $[ ]_N$

features ( $N, D$ )

$\downarrow$   
 $W\text{-proj } (D, H)$   
 $h_0 (N, H)$

The value of  
all the words

soft.-start

$\rightarrow$  start ( $N, 1$ )

$\downarrow$  word\_embedding\_forward ( $x, W$ )

$x_0 (N, D)$

( $N, T$ )

( $V, D$ )

$h_0 \quad x_0$

$\downarrow$   
rnn-step-forward ( $x, \text{prev-h}, \dots$ )

$\downarrow$   
next-h ( $N, H$ )

affine\_forward ( $\text{next-h}, W\text{-vocab}, b\text{-vocab}$ )

$\xrightarrow{\text{argmax}}$  out ( $N, V$ )

$x (N, 1)$

## LSTM - Captioning

$$a = W_x x_t + W_h h_{t-1} + b$$

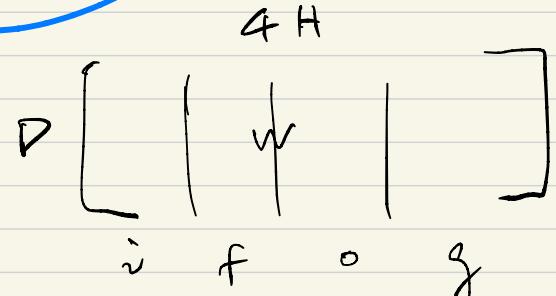
$(N, D)$        $(D, 4H)$        $(4H, 4H)$

$\downarrow$        $\downarrow$        $\downarrow$

$(D, 4H)$        $(N, H)$

$(N, 4H)$

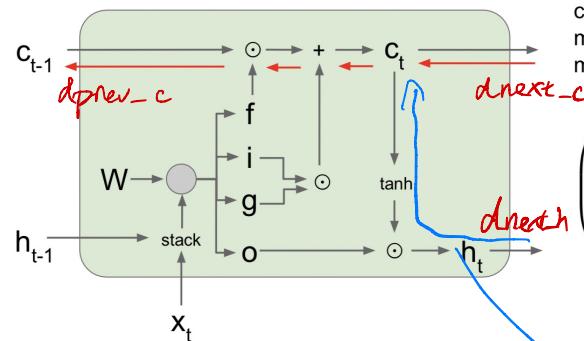
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \circ W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$



$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

Long Short Term Memory (LSTM): Gradient Flow  
[Hochreiter et al., 1997]



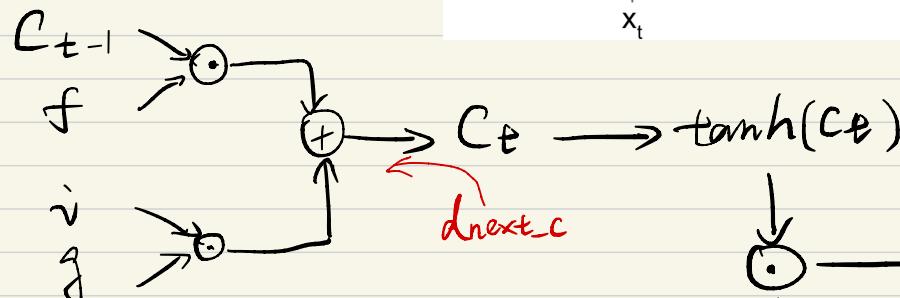
Backpropagation from  $c_t$  to  $c_{t-1}$ : only elementwise multiplication by  $f$ , no matrix multiply by  $W$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

Backward pass



If it's necessary?

$O$   
 $(N, H)$

$x \rightarrow (N, D)$

Input  $\rightarrow d_{\text{next\_h}}, d_{\text{next\_c}}$

$$[x] \cdot [N, D] \quad [N, 4H]$$

$$d_0 = d_{\text{next\_h}} \odot \tanh(c_t)$$

$$dW_x = x^T [d_0 \odot (0 \odot (1 - 0))] \quad (D, N) \quad (N, A) \quad (3)$$

$$df = d_{\text{next\_c\_all}} \odot p_{\text{prev\_c}}$$

$$dW_f = x^T [df \odot f \odot (1 - f)] \quad (2)$$

$$di = d_{\text{next\_c\_all}} \odot g$$

$$dW_{xi} = x^T [di \odot i \odot (1 - i)] \quad (1)$$

$$dg = d_{\text{next\_c\_all}} \odot i \quad \tanh(xW_{xg})$$

$$dW_{xg} = x^T [dg \odot (1 - g^2)] \quad (4)$$

这里需要考虑  $d_{\text{next\_h}}$ ? YES

$\tanh'(c_{\text{next\_c}})$

$$d_{\text{next\_c\_all}} = d_{\text{next\_c}} +$$

$$d_{\text{next\_h}} \odot \odot (1 - \tanh^2(c_{\text{next\_c}}))$$

$$d_{\text{prev\_c}} = d_{\text{next\_c\_all}} \odot f$$

(1)(2)(3)(4) 分别取  $x^T$  中得到  $dW_x$ . 以  $w$  concatenate

起来一起下来

$$dW_x = x^T \cdot \quad (D, N)$$



$$dW_h = \text{prev\_h}^T \cdot d\text{gate}$$

$$dx = d\text{gate} \cdot W_x^T$$

$$d\text{prev\_h} = d\text{gate} \cdot W_h^T \quad (N, 4H)$$

$$db = \text{np.sum}(d\text{gate}, \text{axis}=0)$$

$$Y = xw$$

$$dY = dw$$

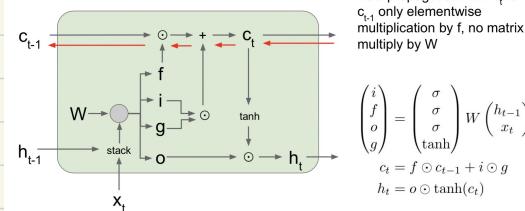
$$\frac{\partial Y}{\partial w} = (x)^T$$

$$dY = dx \cdot w$$

$$\text{tr}(dY) = \text{tr}(w^T dx)$$

$$\frac{\partial Y}{\partial x^T} = w \Rightarrow \frac{\partial Y}{\partial x} = w$$

Long Short Term Memory (LSTM): Gradient Flow  
[Hochreiter et al., 1997]



Backpropagation from  $c_t$  to  $c_{t-1}$  only elementwise multiplication by  $f$ , no matrix multiplication by  $W$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

# Network Visualization

1. **Saliency Maps:** Saliency maps are a quick way to tell which part of the image influenced the classification decision made by the network.
2. **Fooling Images:** We can perturb an input image so that it appears the same to humans, but will be misclassified by the pretrained network.
3. **Class Visualization:** We can synthesize an image to maximize the classification score of a particular class; this can give us some sense of what the network is looking for when it classifies images of that class.

— Create a session that dynamically allocates memory.  
↳ allowing gpu memory growth.

默认情况下，TensorFlow 会映射进程可见的所有 GPU 的几乎所有 GPU 内存（取决于 CUDA\_VISIBLE\_DEVICES）。通过减少 内存碎片，可以更有效地使用设备上相对宝贵的 GPU 内存资源。

```
config = tf.ConfigProto()  
config.gpu_options.allow_growth = True  
session = tf.Session(config=config, ...)
```

↳ memory fragmentation.  
a phenomenon in which storage space is used inefficiently, reducing capacity or performance and often both.

— `tf.get_variable(name, shape=[..], ...)`

创建一个 variable。默认的 dtype 为 `tf.float32`，初始化通过 `tf.glorot_uniform_initializer` 随机得到。

`tf.variable_scope(name, reuse=...)` 将 `tf.Variable` 对象隐式地封装。即在不同的 scope 下可选择重新使用现有变量还是创建新变量。

```
def my_image_filter(input_images):  
    with tf.variable_scope("conv1"):  
        # Variables created here will be named "conv1/weights", "conv1/biases".  
        relu1 = conv_relu(input_images, [5, 5, 32, 32], [32])  
    with tf.variable_scope("conv2"):  
        # Variables created here will be named "conv2/weights", "conv2/biases".  
        return conv_relu(relu1, [5, 5, 32, 32], [32])
```

`tf.Variable`

这里“conv\_relu”为创建一个卷积层的函数，其中包括创建 weights 和 biases。由于在不同的层时都需要创建新的变量，于是用不同的 scope 隐式地进行封装。

如果要共享变量，则可以将 “reuse=True” 并将其 scope 名相同，或调用 `scope.reuse_variables()`

```

with tf.variable_scope("model"):
    output1 = my_image_filter(input1)
with tf.variable_scope("model", reuse=True):
    output2 = my_image_filter(input2)

```

```

with tf.variable_scope("model") as scope:
    output1 = my_image_filter(input1)
    scope.reuse_variables()
    output2 = my_image_filter(input2)

```

## Saliency Maps

$\frac{\partial \text{score}}{\partial \text{pixels of image}}$

Compute the gradient of the unnormalized score corresponding to the correct class w.r.t the pixels of the image.

image ( $l, w, 3$ ) gradient ( $l, w, 3$ )  $\xrightarrow{\max_{\text{axis}=2}}$  saliency map ( $l, w$ )

## Class Visualization

By starting with a random noise image and performing gradient ascent on a target class, we can generate an image that the network will recognize as the target class.

$$\text{regularizer, } R(I) = \lambda \|I\|_2^2$$

$$I^* = \arg \max_I (\underbrace{S_y(I)}_{\text{Score that a CNN assigns to image } I \text{ for class } y} - \underbrace{R(I)}_{\text{regularizer}})$$

$I^*$  generated image

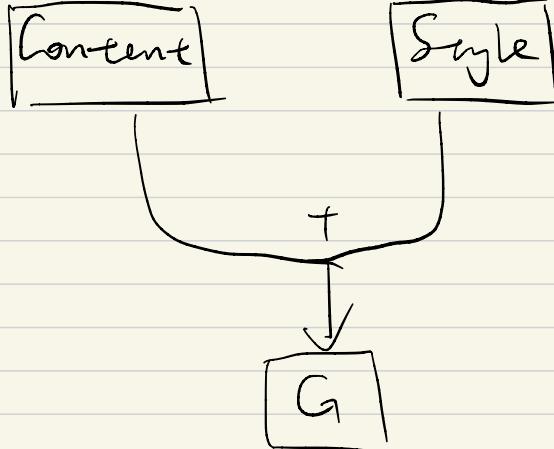
Score that a CNN assigns to image  $I$  for class  $y$  (raw unnormalized scores)

## Fooling Image

Given an image and a target class, perform gradient ascent over the image to maximize the target class, stopping when the network classifies the image as the target class.



## Style Transfer



Cost function:

$$L(G) = \alpha L_{\text{content}}(C, G) + \beta L_{\text{style}}(S, G)$$

Generated image

Find the generated image  $G$ :

- ① Initialize  $G$  randomly
- ② Use gradient descent to minimize  $L(G)$   
update the pixels in  $G$

$$G := G - \frac{\partial}{\partial G} L(G)$$

## Content loss

Content loss measures how much the generated image differs from the feature maps of the source image.

in hidden layer  $L$ , feature maps  $A^L \in \mathbb{R}^{1 \times H_L \times W_L \times C_L}$

reshape  $\rightarrow$  feature map for current image  $F^L \in \mathbb{R}^{M_L \times C_L}$   
feature map for content source  $P^L \in \mathbb{R}^{M_L \times C_L}$

$$M_L = H_L \times W_L$$

$M_L \begin{bmatrix} | & | & | & | \\ | & | & | & | \end{bmatrix}$  feature content loss of 不是子系统  $\rightarrow$  feature map reshape  $\hat{x} = \hat{x}^L$ ,  $\hat{x}^L$  里面包含 style loss 方便

Each column of  $F^L$  or  $P^L$  represents the vectorized activations of a particular filter, convolved over all positions of the image.

$$L_{\text{content}}(c, G) = \|F^L - P^L\|_2^2 = \sum_{i,j} (F_{ij}^L - P_{ij}^L)^2$$

## Syle loss

Define style as correlation between activations across channels.

Aside:

### Matrix Norm

① Matrix norms induced by vector norms.

(induced norm)

$$\|A\|_p = \sup_{x \neq 0} \frac{\|Ax\|_p}{\|x\|_p}$$

这里是常数的向量的  $L_p$  范数

$A \rightarrow \mathbb{R}^{m \times n}$

$x \rightarrow \mathbb{R}^n$

spectral cases:

$$\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}| \rightarrow \text{maximum column sum}$$

(sum across rows)

$$\|A\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}| \rightarrow \text{maximum row sum}$$

(sum across columns)

$$\|x\|_1 = \sum_i |x_i|$$

$$\|x\|_\infty = \max_i |x_i|$$

$$\|A\|_2 = \sigma_{\max}(A) = \sqrt{\lambda_{\max}(A^T A)}$$

largest singular value of matrix  $A$ .

spectral norm (the operator norm induced by the vector 2-norm)

## ② "Entrywise" matrix norms

treat an  $m \times n$  matrix as a vector of size  $m \cdot n$ , and use vector norms.

$$\|A\|_p = \|\text{vec}(A)\|_p = \left( \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^p \right)^{\frac{1}{p}}$$

$$\|A\|_{p,q} = \left[ \sum_{j=1}^n \left( \sum_{i=1}^m |a_{ij}|^p \right)^{\frac{q}{p}} \right]^{\frac{1}{q}}$$

Spectral cases:

$$\begin{aligned} - \text{Frobenius norm} \rightarrow \|A\|_F &= \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} = \sqrt{\text{trace}(A^* A)} \\ &= \sqrt{\sum_{i=1}^{\min\{m, n\}} \sigma_i^2(A)} \end{aligned}$$

$\sigma_i(A)$  are the singular values of  $A$ .

induced norm

$$\|A\|_2 = \sigma_{\max}(A) \leq \|A\|_F \quad (\text{Equation holds only if } A \text{ is a rank-one matrix or a zero matrix})$$

$\|A\|_2$  is Frobenius norm for  $p=2$ .

$\|A\|_2$  is induced norm  $\Leftrightarrow \|A\|_F$

$$-\text{Max norm} \rightarrow \|A\|_{\max} = \max_{i,j} |a_{ij}|$$

## ③ Schatten norms

The Schatten  $p$ -norms arise when applying the  $p$ -norm to the vector of singular values of a matrix.

$$\|A\|_p = \left( \sum_{i=1}^{\min\{m, n\}} \sigma_i^p(A) \right)^{\frac{1}{p}}$$

Spectral cases:

$p=1 \rightarrow$  nuclear norm (trace norm)

$p=2 \rightarrow$  Frobenius norm

$p=\infty \rightarrow$  Spectral norm

## Vector Norm

$$\|x\|_p = \left( \sum_i |x_i|^p \right)^{\frac{1}{p}}$$

$$\|x\|_1 = \sum_i |x_i| \rightarrow L^1 \text{ 范数. 当机器学习问题中零和非零元素之间差异非常大时会使用 } L^1 \text{ 范数}$$

$$\|x\|_\infty = \max_i |x_i|$$

$$\|x\|_2 = \sqrt{\sum_i |x_i|^2} \rightarrow L^2 \text{ 范数 (Euclidean Norm)}$$

$$\|x\|_2^2 = \sum_i |x_i|^2 = x^T x$$

$L^2$  范数经常用来衡量向量的大小，但在原点附近增长缓慢，因为平方  $L^2$  范数对  $x$  中每个元素的导数只取决于对应元素本身。

前面提到了 Syle 定义为 "correlations between activations across channels" 因为不同 channels 对应着不同的 filters.

这里在实现应用中使用了 Gram matrix 作为 covariance matrix 的近似。我们希望的是 generated image 的 activations 不同 channels 之间的相关性  $\mu_x$  和 Syle image 的相关性。

Covariance matrix  $\rightarrow \text{Cov}[x, x]$

$$= E[(x - \mu_x)(x - \mu_x)^T]$$

$$= \underline{E[x x^T]} - \mu_x \mu_x^T$$

Gram matrix  $\rightarrow x x^T$  (这里 in matrix  $x$  定义为

$$\begin{bmatrix} -x_1^T- \\ -x_2^T- \\ \vdots \\ -x_n^T- \end{bmatrix})$$

对于向量来说 Gram matrix 就是内积

由于 assignment 中定义的  $\|\cdot\|_C$ ，所以这里：

$$G_g^L = F^T F$$

$$G_s^L = S^T S$$

$$\Delta_{Syle}(S, G) = \|G_g^L - G_s^L\|_2^2 = \sum_{ij} (G_{g_{ij}}^L - G_{s_{ij}}^L)^2$$

使用 channels 用 in 预先令风格为 style 而该是标准的解译为，而表示那些高层的纹理或元素倾向于同时或不同的出现的程度。

## Total-variation Regularization.

It's also helpful to encourage "smoothness" in the images.

但这是否在 loss function 中加入对 wiggles 的 "total variation" in penalty

total-variation



sum of the squares of differences in the pixel values for all pairs of pixels that are next to each other (horizontally or vertically).

img: (1, H, W, 3)

$$\begin{aligned} & \underset{\substack{\text{O}^{(i, j-1, c)} \\ \downarrow \\ \text{O} \leftarrow \text{O}^{(i, j, c)}}}{\text{current pixel}} = \text{tf.reduce\_sum}((\text{img}[:, :, :, :] - \text{img}[:, :, :-1, :])^2) \\ & + \text{tf.reduce\_sum}((\text{img}[:, :, :, :] - \text{img}[:, :, :, :-1])^2) \end{aligned}$$

## Generative Adversarial Networks (GANs)

Aside: 交叉熵和最大似然估计。

### ① 最大似然估计

一组会有 m 个样本的数据集  $X = \{x^{(1)}, \dots, x^{(m)}\}$ , 独立地由未知的真分布  $P_{data}(x)$  生成。

令  $P_{model}(x; \theta)$  是由  $\theta$  确定在相同空间上的概率分布。 $P_{model}(x; \theta)$  将任意输入  $x$  映射到实数来估计真实概率  $P_{data}(x)$ 。

$$\begin{aligned} \text{对 } \theta \text{ 的最大似然估计定义为} \rightarrow \theta_{ML} &= \arg \max_{\theta} P_{model}(X; \theta) \\ &= \arg \max_{\theta} \prod_{i=1}^m P_{model}(x^{(i)}; \theta) \end{aligned}$$

最大似然估计是找到使得模型而产生出实际样本概率的而最大化模型参数的方法。

为了简化计算, 同时保证数值不会下溢, 将乘积转化为对数求和的形式。

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log P_{\text{model}}(x^{(i)}; \theta)$$

因为当重新编写代码时， $\operatorname{argmax}$  不变，所以可以将  $m$  提到从训练集中经验分布  $\hat{P}_{\text{data}}$  相关的期望作为准则：

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} E_{x \sim \hat{P}_{\text{data}}} \log P_{\text{model}}(x; \theta)$$

Population  $\rightarrow$  Observation  $\rightarrow$  model

$$P_{\text{data}}(x) \quad \hat{P}_{\text{data}}(x) \quad P_{\text{model}}(x; \theta)$$

从 population 到 observation 再到 model 是怎样的关系？

最大似然估计的第一步就是将最大化训练集上经验分布  $\hat{P}_{\text{data}}$   
和模型分布之间的差异。

这种差异可以通过 KL 散度度量：

$$D_{KL}(\hat{P}_{\text{data}} \parallel P_{\text{model}}) = E_{x \sim \hat{P}_{\text{data}}} [\log \hat{P}_{\text{data}}(x) - \log P_{\text{model}}(x)]$$

$$\min D_{KL}(\hat{P}_{\text{data}} \parallel P_{\text{model}}) = \min -E_{x \sim \hat{P}_{\text{data}}} [\log P_{\text{model}}(x)]$$

最大化 KL 故障其实就是在最小化分布之间的交叉熵。任何一个由负对数似然函数的损失都是定义在训练集经验分布和定义在模型上的平均分布之间的交叉熵。

## (2) 交叉熵

自信息  $I(x) = -\log P(x) \rightarrow$  事件  $x$  的信息量

$$\text{香农熵: } H(x) = E_{x \sim P}[I(x)] = -E_{x \sim P}[\log P(x)]$$

$\hookrightarrow$  真值分布  $P$  中事件产生的期望信息量

The Shannon entropy of a distribution is the expected amount of information in an event drawn from that distribution.

$$KL\text{ 故度} : D_{KL}(P||Q) = E_{x \sim P} \left[ \log \frac{P(x)}{Q(x)} \right] = E_{x \sim P} [\log P(x) - \log Q(x)]$$

→ 從這下隨機變量  $x$  的二種分布  $P(x)$  和  $Q(x)$  之差

$$\text{交叉熵} : H(P, Q) = H(P) + D_{KL}(P||Q)$$

$$= -E_{x \sim P} [\log Q(x)] = -\sum_x P(x) \log Q(x)$$

針對  $Q$  小化交叉熵等價于最小化 KL 故度。

