

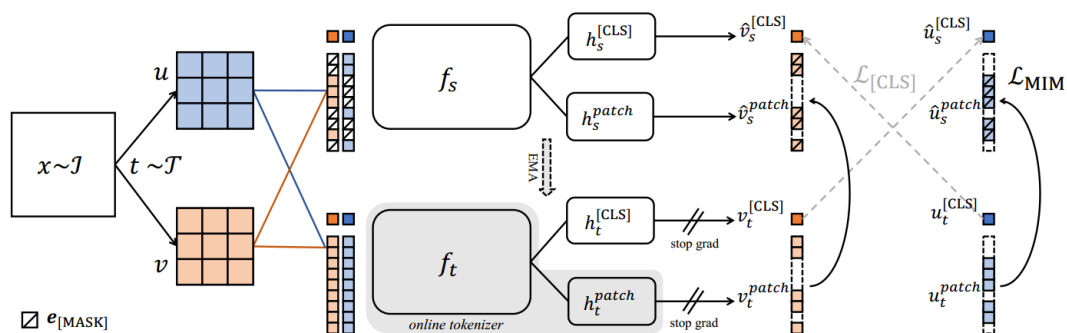
iBOT论文阅读

文章标题: *iBOT: Image Bert pre-training with Online Tokenizer*, 可以理解为使用在线tokenizer进行图像BERT式预训练

写在前面:

iBOT可以视为BEIT和DINO的结合, 在iBOT中学生网络和教师网络的输出会有两部分, 一部分是和DINO中类似的 $h^{[CLS]}$ 一部分是和BEIT中类似的 h^{patch}

iBOT的模型结构图和算法流程:



Algorithm 1: iBOT PyTorch-like Pseudocode w/o multi-crop augmentation

Input:

```
 $g_s, g_t$  ; // student and teacher network  
 $C, C'$  ; // center on [CLS] token and patch tokens  
 $\tau_s, \tau_t$  ; // temperature on [CLS] token for student and teacher network  
 $\tau'_s, \tau'_t$  ; // temperature on patch tokens for student and teacher network  
 $l$  ; // momentum rate for network  
 $m, m'$  ; // momentum rates for center on [CLS] token and patch tokens
```

$g_t.params = g_s.params$

for x in loader **do**

```
 $u, v = \text{augment}(x), \text{augment}(x)$  ; // random views  
 $\hat{u}, m_u = \text{blockwise\_mask}(u)$  ; // random block-wise masking  
 $\hat{v}, m_v = \text{blockwise\_mask}(v)$  ; // random block-wise masking  
  
 $\hat{u}_s^{[CLS]}, \hat{u}_s^{\text{patch}} = g_s(\hat{u}, \text{return\_all\_tok}=\text{true})$  ; //  $[n, K], [n, S^2, K]$   
 $\hat{v}_s^{[CLS]}, \hat{v}_s^{\text{patch}} = g_s(\hat{v}, \text{return\_all\_tok}=\text{true})$  ; //  $[n, K], [n, S^2, K]$   
  
 $u_t^{[CLS]}, u_t^{\text{patch}} = g_t(u, \text{return\_all\_tok}=\text{true})$  ; //  $[n, K], [n, S^2, K]$   
 $v_t^{[CLS]}, v_t^{\text{patch}} = g_t(v, \text{return\_all\_tok}=\text{true})$  ; //  $[n, K], [n, S^2, K]$   
  
 $\mathcal{L}_{[CLS]} = H(\hat{u}_s^{[CLS]}, v_t^{[CLS]}, C, \tau_s, \tau_t) / 2 + H(\hat{v}_s^{[CLS]}, u_t^{[CLS]}, C, \tau_s, \tau_t) / 2$   
 $\mathcal{L}_{\text{MIM}} = (m_u \cdot H(\hat{u}_s^{\text{patch}}, u_t^{\text{patch}}, C', \tau'_s, \tau'_t).sum(dim=1) / m_u.sum(dim=1) / 2$   
           $+ (m_v \cdot H(\hat{v}_s^{\text{patch}}, v_t^{\text{patch}}, C', \tau'_s, \tau'_t).sum(dim=1) / m_v.sum(dim=1) / 2$   
           $(\mathcal{L}_{[CLS]}.mean() + \mathcal{L}_{\text{MIM}}.mean()).backward()$   
  
update( $g_s$ ) ; // student, teacher and center update  
 $g_t.params = l \cdot g_t.params + (1 - l) \cdot g_s.params$   
 $C = m \cdot C + (1 - m) \cdot \text{cat}([u_t^{[CLS]}, v_t^{[CLS]}]).mean(dim=0)$   
 $C' = m' \cdot C' + (1 - m') \cdot \text{cat}([u_t^{\text{patch}}, v_t^{\text{patch}}]).mean(dim=(0, 1))$ 
```

end**def** $H(s, t, c, \tau_s, \tau_t)$:

```
 $t = t.detach()$  ; // stop gradient  
 $s = \text{softmax}(s / \tau_s, dim=1)$   
 $t = \text{softmax}((t - c) / \tau_t, dim=1)$  ; // center + sharpen  
return  $-(t \cdot \log(s)).sum(dim=-1)$ 
```

算法流程:

1. 与DINO类似，一开始学生网络和教师网络的参数是一致的
2. 与DINO类似，每个图像 x ，会经过两种不同的增广方式得到 u 和 v
3. 与DINO不同的是，进入学生网络的图像 \hat{u} 和 \hat{v} 是 u 和 v 被mask掉的一部分后的图像
4. 在iBOT中学生网络和教师网络的输出会有两部分：
 - 第一部分与DINO类似（暂时理解为图像的全局特征，对图像整体的编码），学生网络输出 $\hat{u}_s^{[CLS]}$ 和 $\hat{v}_s^{[CLS]}$
教师网络输出 $u_t^{[CLS]}$ 和 $v_t^{[CLS]}$
 - 第二部分与BEiT类似（被mask掉的patch特征），学生网络输出 \hat{u}_s^{patch} 和 \hat{v}_s^{patch}
教师网络输出 u_t^{patch} 和 v_t^{patch}

第一部分的损失函数:

$$\mathcal{L}_{[CLS]} = H(\hat{u}_s^{[CLS]}, v_t^{[CLS]}, C, \tau_s, \tau_t) / 2 + H(\hat{v}_s^{[CLS]}, u_t^{[CLS]}, C, \tau_s, \tau_t) / 2$$

和DINO文章中的损失函数是一致的：

```
x1, x2 = augment(x), augment(x) # random views

s1, s2 = gs(x1), gs(x2) # student output n-by-K
t1, t2 = gt(x1), gt(x2) # teacher output n-by-K

loss = H(t1, s2)/2 + H(t2, s1)/2
```

H 表示交叉熵损失函数，更准确的说是带有温度参数（锐度）和偏置量（居中）的交叉熵损失函数。位置和居中的平衡避免了模型在训练过程中崩溃。

第二部分的损失函数：

$$\mathcal{L}_{MIM} = (m_u \cdot H(\hat{u}_s^{\text{patch}}, u_t^{\text{patch}}, C', \tau'_s, \tau'_t). \text{sum}(\text{dim}=1) / m_u. \text{sum}(\text{dim}=1) / 2 \\ + (m_v \cdot H(\hat{v}_s^{\text{patch}}, v_t^{\text{patch}}, C', \tau'_s, \tau'_t). \text{sum}(\text{dim}=1) / m_v. \text{sum}(\text{dim}=1) / 2$$

计算被mask掉patch的特征相似度（ m_u 和 m_v 的参数含义可能是被mask掉的比例？这里需要再了解一下）

总体损失函数：

$$L_{[CLS]}. \text{mean}() + L_{MIM}. \text{mean}()$$

5. 参数更新：

- 学生网络的参数更新是根据损失函数反向传播结果优化的
- 教师网络参数更新与DINO类似，采用动量法，结合自己本来的参数以及学生网络的参数
- 同样的，偏置量 C 和 C' 也是采用动量法更新

采用动量法更新的公式如下图所示：

$$g_t. \text{params} = l \cdot g_t. \text{params} + (1 - l) \cdot g_s. \text{params} \\ C = m \cdot C + (1 - m) \cdot \text{cat}([u_t^{[CLS]}, v_t^{[CLS]}]). \text{mean}(\text{dim}=0) \\ C' = m' \cdot C' + (1 - m') \cdot \text{cat}([u_t^{\text{patch}}, v_t^{\text{patch}}]). \text{mean}(\text{dim}=(0, 1))$$