

用Pytorch轻松实现28个视觉Transformer，开源库 timm 了解一下！（附代码解读）

原创

CV开发者都爱看的

极市平台

2021-02-19 22:00:00

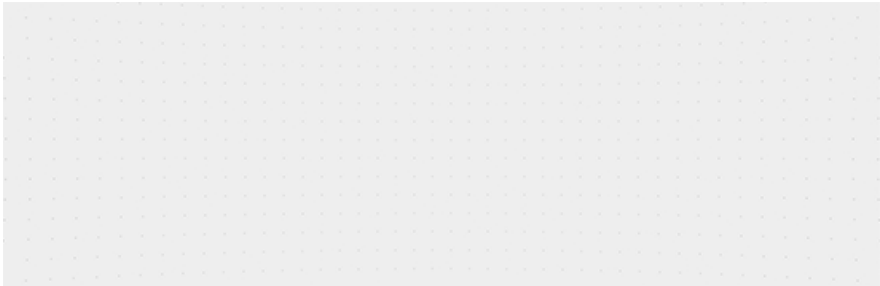
手机阅读

𐄞

收录于话题

#Transformer

↑ 点击蓝字 关注极市平台



作者 | 科技猛兽

审稿 | 邓富城

编辑 | 极市平台

极市导读

本文将介绍一个优秀的PyTorch开源库——timm库，并对其中的vision transformer.py代码进行了详细解读。 >>加入极市CV技术交流群，走在计算机视觉的最前沿

Transformer 架构早已在自然语言处理任务中得到广泛应用，但在计算机视觉领域中仍然受到限制。在计算机视觉领域，目前已有大量工作表明模型对 CNN 的依赖不是必需的，当直接应用于图像块序列时，Transformer 也能很好地执行图像分类任务。

本文将简要介绍了优秀的 PyTorch Image Model 库：timm库。与此同时，将会为大家详细介绍其中的视觉Transformer代码以及一个优秀的视觉Transformer 的PyTorch实现，以帮助大家更快地开展相关实验。

什么是timm库？


PyTorchImageModels，简称timm，是一个巨大的PyTorch代码集合，包括了一系列：

- image models
- layers
- utilities
- optimizers
- schedulers
- data-loaders / augmentations
- training / validation scripts

旨在将各种SOTA模型整合在一起，并具有复现ImageNet训练结果的能力。

timm库作者是来自加拿大温哥华的Ross Wightman。

 壹伴图


 极市平台
extreme


月发文数目： **


月平均阅读： **

文章工具

已发

 采集图文

 合成多

 采集样式

 查看样



作者github链接：

<https://github.com/rwightman>

timm库链接：

<https://github.com/rwightman/pytorch-image-models>

所有的PyTorch模型及其对应arxiv链接如下：

- Big Transfer ResNetV2 (BiT) - <https://arxiv.org/abs/1912.11370>
- CspNet (Cross-Stage Partial Networks) - <https://arxiv.org/abs/1911.11929>
- DeiT (Vision Transformer) - <https://arxiv.org/abs/2012.12877>
- DenseNet - <https://arxiv.org/abs/1608.06993>
- DLA - <https://arxiv.org/abs/1707.06484>
- DPN (Dual-Path Network) - <https://arxiv.org/abs/1707.01629>

timm库特点

所有的模型都有默认的API：

- accessing/changing the classifier - `get_classifier` and `reset_classifier`
- 只对features做前向传播 - `forward_features`

所有模型都支持多尺度特征提取 (feature pyramids) (通过create_model函数)：

- `create_model(name, features_only=True, out_indices=..., output_stride=...)`

`out_indices` 指定返回哪个feature maps to return, 从0开始, `out_indices[i]` 对应着 $C(i + 1)$ feature level。

`output_stride` 通过dilated convolutions控制网络的output stride。大多数网络默认 stride 32。

所有的模型都有一致的pretrained weight loader, adapts last linear if necessary。

训练方式支持：

- NVIDIA DDP w/ a single GPU per process, multiple processes with APEX present (AMP mixed-precision optional)
- PyTorch DistributedDataParallel w/ multi-gpu, single process (AMP disabled as it crashes when enabled)
- PyTorch w/ single GPU single process (AMP optional)

动态的全局池化方式可以选择： average pooling, max pooling, average + max, or concat([average, max]), 默认是adaptive average。

Schedulers:

Schedulers 包括 [step](#) , [cosine](#) w/ restarts, [tanh](#) w/ restarts, [plateau](#) 。

Optimizer:

- [rmsprop_tf](#) adapted from PyTorch RMSProp by myself. Reproduces much improved Tensorflow RMSProp behaviour.
- [radam](#) by Liyuan Liu (<https://arxiv.org/abs/1908.03265>)
- [novograd](#) by Masashi Kimura (<https://arxiv.org/abs/1905.11286>)
- [lookahead](#) adapted from impl by Liam (<https://arxiv.org/abs/1907.08610>)
- [fused<name>](#) optimizers by name with NVIDIA Apex installed
- [adamp](#) and [sgdp](#) by Naver ClovAI (<https://arxiv.org/abs/2006.08217>)
- [adafactor](#) adapted from FAIRSeq impl (<https://arxiv.org/abs/1804.04235>)
- [adahessian](#) by David Samuel (<https://arxiv.org/abs/2006.00719>)

timm库 vision_transformer.py代码解读

代码来自：

https://github.com/rwightman/pytorch-image-models/blob/master/timm/models/vision_transformer.py

对应的论文是ViT，是除了官方开源的代码之外的又一个优秀的PyTorch implement。

An Image Is Worth 16 x 16 Words: Transformers for Image Recognition at Scale

An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale

<https://arxiv.org/abs/2010.11929>

另一篇工作DeiT也大量借鉴了timm库这份代码的实现：

Training data-efficient image transformers & distillation through attention

Training data-efficient image transformers & distillation through attention

<https://arxiv.org/abs/2012.12877>

[vision_transformer.py](#):

代码中定义的变量的含义如下：

img_size: tuple 类型，里面是int类型，代表输入的图片大小，默认是 **224**。

patch_size: tuple 类型，里面是int类型，代表Patch的大小，默认是 **16**。

in_chans: int 类型，代表输入图片的channel数，默认是**3**。

num_classes: int 类型classification head的分类数，比如CIFAR100就是100，默认是 **1000**。

embed_dim: int 类型Transformer的embedding dimension，默认是 **768**。

depth: int 类型，Transformer的Block的数量，默认是 **12**。

num_heads: int 类型，attention heads的数量，默认是**12**。

mlp_ratio: int 类型，mlp hidden dim/embedding dim的值，默认是 **4**。

qkv_bias: bool 类型，attention模块计算qkv时需要bias吗，默认是 **True**。

qk_scale: 一般设置成 **None** 就行。

drop_rate: float 类型，dropout rate，默认是 **0**。

attn_drop_rate: float 类型，attention模块的dropout rate，默认是 **0**。

drop_path_rate: float 类型，默认是 **0**。

hybrid_backbone: nn.Module 类型，在把图片转换成Patch之前，需要先通过一个Backbone吗？默认是 **None**。

如果是None，就直接把图片转化成Patch。

如果不是None，就先通过这个Backbone，再转化成Patch。

norm_layer: nn.Module 类型，归一化层类型，默认是 **None**。

1. 导入必要的库和模型：

```
1 import math
2 import logging
3 from functools import partial
4 from collections import OrderedDict
5
6 import torch
7 import torch.nn as nn
8 import torch.nn.functional as F
9
10 from timm.data import IMAGENET_DEFAULT_MEAN, IMAGENET_DEFAULT_STD
11 from .helpers import load_pretrained
12 from .layers import StdConv2dSame, DropPath, to_2tuple, trunc_normal_
13 from .resnet import resnet26d, resnet50d
14 from .resnetv2 import ResNetV2
15 from .registry import register_model
```

2. 定义一个字典，代表标准的模型，如果需要更改模型超参数只需要改变_cfg

的传入的参数即可。

```
1 def _cfg(url='', **kwargs):
```

```

2     return {
3         'url': url,
4         'num_classes': 1000, 'input_size': (3, 224, 224), 'pool_size': None,
5         'crop_pct': .9, 'interpolation': 'bicubic',
6         'mean': IMAGENET_DEFAULT_MEAN, 'std': IMAGENET_DEFAULT_STD,
7         'first_conv': 'patch_embed.proj', 'classifier': 'head',
8         **kwargs
9     }

```

3. default_cfgs代表支持的所有模型，也定义成字典的形式：

vit_small_patch16_224里面的small代表小模型。

ViT的第一步要把图片分成一个个**patch**，然后把这些patch组合在一起作为对图像的序列化操作，比如一张224 × 224的图片分成大小为16 × 16的patch，那一共可以分成196个。所以这个图片就序列化成了(196, 256)的tensor。所以这里的：

16：就代表patch的大小。

224：就代表输入图片的大小。

按照这个命名方式，支持的模型有：vit_base_patch16_224，vit_base_patch16_384等等。

后面的vit_deit_base_patch16_224等等模型代表DeiT这篇论文模型。

```

1  default_cfgs = {
2      # patch models (my experiments)
3      'vit_small_patch16_224': _cfg(
4          url='https://github.com/rwightman/pytorch-image-models/releases/download/v0.1-weights/vit_small_patch16_224.pth',
5      ),
6
7      # patch models (weights ported from official Google JAX impl)
8      'vit_base_patch16_224': _cfg(
9          url='https://github.com/rwightman/pytorch-image-models/releases/download/v0.1-vitjx/jx_vit_base_patch16_224_in21k.pth',
10         mean=(0.5, 0.5, 0.5), std=(0.5, 0.5, 0.5),
11     ),
12     'vit_base_patch32_224': _cfg(
13         url='', # no official model weights for this combo, only for in21k
14         mean=(0.5, 0.5, 0.5), std=(0.5, 0.5, 0.5)),
15     'vit_base_patch16_384': _cfg(
16         url='https://github.com/rwightman/pytorch-image-models/releases/download/v0.1-vitjx/jx_vit_base_patch16_384_in21k.pth',
17         input_size=(3, 384, 384), mean=(0.5, 0.5, 0.5), std=(0.5, 0.5, 0.5), crop_pct=1.0),
18     'vit_base_patch32_384': _cfg(
19         url='https://github.com/rwightman/pytorch-image-models/releases/download/v0.1-vitjx/jx_vit_base_patch32_384_in21k.pth',
20         input_size=(3, 384, 384), mean=(0.5, 0.5, 0.5), std=(0.5, 0.5, 0.5), crop_pct=1.0),
21     'vit_large_patch16_224': _cfg(
22         url='https://github.com/rwightman/pytorch-image-models/releases/download/v0.1-vitjx/jx_vit_large_patch16_224_in21k.pth',
23         mean=(0.5, 0.5, 0.5), std=(0.5, 0.5, 0.5)),
24     'vit_large_patch32_224': _cfg(
25         url='', # no official model weights for this combo, only for in21k
26         mean=(0.5, 0.5, 0.5), std=(0.5, 0.5, 0.5)),
27     'vit_large_patch16_384': _cfg(
28         url='https://github.com/rwightman/pytorch-image-models/releases/download/v0.1-vitjx/jx_vit_large_patch16_384_in21k.pth',
29         input_size=(3, 384, 384), mean=(0.5, 0.5, 0.5), std=(0.5, 0.5, 0.5), crop_pct=1.0),
30     'vit_large_patch32_384': _cfg(
31         url='https://github.com/rwightman/pytorch-image-models/releases/download/v0.1-vitjx/jx_vit_large_patch32_384_in21k.pth',
32         input_size=(3, 384, 384), mean=(0.5, 0.5, 0.5), std=(0.5, 0.5, 0.5), crop_pct=1.0),

```

```

33
34 # patch models, imagenet21k (weights ported from official Google JAX impl)
35 'vit_base_patch16_224_in21k': _cfg(
36     url='https://github.com/rwightman/pytorch-image-models/releases/download/v0.1-vitjx/jx_vit_base
37     num_classes=21843, mean=(0.5, 0.5, 0.5), std=(0.5, 0.5, 0.5)),
38 'vit_base_patch32_224_in21k': _cfg(
39     url='https://github.com/rwightman/pytorch-image-models/releases/download/v0.1-vitjx/jx_vit_base
40     num_classes=21843, mean=(0.5, 0.5, 0.5), std=(0.5, 0.5, 0.5)),
41 'vit_large_patch16_224_in21k': _cfg(
42     url='https://github.com/rwightman/pytorch-image-models/releases/download/v0.1-vitjx/jx_vit_lar
43     num_classes=21843, mean=(0.5, 0.5, 0.5), std=(0.5, 0.5, 0.5)),
44 'vit_large_patch32_224_in21k': _cfg(
45     url='https://github.com/rwightman/pytorch-image-models/releases/download/v0.1-vitjx/jx_vit_lar
46     num_classes=21843, mean=(0.5, 0.5, 0.5), std=(0.5, 0.5, 0.5)),
47 'vit_huge_patch14_224_in21k': _cfg(
48     url='', # FIXME I have weights for this but > 2GB limit for github release binaries
49     num_classes=21843, mean=(0.5, 0.5, 0.5), std=(0.5, 0.5, 0.5)),
50
51 # hybrid models (weights ported from official Google JAX impl)
52 'vit_base_resnet50_224_in21k': _cfg(
53     url='https://github.com/rwightman/pytorch-image-models/releases/download/v0.1-vitjx/jx_vit_base
54     num_classes=21843, mean=(0.5, 0.5, 0.5), std=(0.5, 0.5, 0.5), crop_pct=0.9, first_conv='patch_
55 'vit_base_resnet50_384': _cfg(
56     url='https://github.com/rwightman/pytorch-image-models/releases/download/v0.1-vitjx/jx_vit_base
57     input_size=(3, 384, 384), mean=(0.5, 0.5, 0.5), std=(0.5, 0.5, 0.5), crop_pct=1.0, first_conv=
58
59 # hybrid models (my experiments)
60 'vit_small_resnet26d_224': _cfg(),
61 'vit_small_resnet50d_s3_224': _cfg(),
62 'vit_base_resnet26d_224': _cfg(),
63 'vit_base_resnet50d_224': _cfg(),
64
65 # deit models (FB weights)
66 'vit_deit_tiny_patch16_224': _cfg(
67     url='https://dl.fbaipublicfiles.com/deit/deit_tiny_patch16_224-a1311bcf.pth'),
68 'vit_deit_small_patch16_224': _cfg(
69     url='https://dl.fbaipublicfiles.com/deit/deit_small_patch16_224-cd65a155.pth'),
70 'vit_deit_base_patch16_224': _cfg(
71     url='https://dl.fbaipublicfiles.com/deit/deit_base_patch16_224-b5f2ef4d.pth'),
72 'vit_deit_base_patch16_384': _cfg(
73     url='https://dl.fbaipublicfiles.com/deit/deit_base_patch16_384-8de9b5d1.pth',
74     input_size=(3, 384, 384), crop_pct=1.0),
75 'vit_deit_tiny_distilled_patch16_224': _cfg(
76     url='https://dl.fbaipublicfiles.com/deit/deit_tiny_distilled_patch16_224-b40b3cf7.pth'),
77 'vit_deit_small_distilled_patch16_224': _cfg(
78     url='https://dl.fbaipublicfiles.com/deit/deit_small_distilled_patch16_224-649709d9.pth'),
79 'vit_deit_base_distilled_patch16_224': _cfg(
80     url='https://dl.fbaipublicfiles.com/deit/deit_base_distilled_patch16_224-df68dfff.pth', ),
81 'vit_deit_base_distilled_patch16_384': _cfg(
82     url='https://dl.fbaipublicfiles.com/deit/deit_base_distilled_patch16_384-d0272ac0.pth',
83     input_size=(3, 384, 384), crop_pct=1.0),
84 }

```

4. FFN实现：

```

1 class Mlp(nn.Module):
2     def __init__(self, in_features, hidden_features=None, out_features=None, act_layer=nn.GELU, drop=0):
3         super().__init__()
4         out_features = out_features or in_features
5         hidden_features = hidden_features or in_features
6         self.fc1 = nn.Linear(in_features, hidden_features)
7         self.act = act_layer()
8         self.fc2 = nn.Linear(hidden_features, out_features)
9         self.drop = nn.Dropout(drop)
10
11     def forward(self, x):
12         x = self.fc1(x)
13         x = self.act(x)
14         x = self.drop(x)
15         x = self.fc2(x)
16         x = self.drop(x)
17         return x

```

5. Attention实现：

在python 3.5以后，@是一个操作符，表示矩阵-向量乘法
A@x 就是矩阵-向量乘法A*x: np.dot(A, x)。

```

1 class Attention(nn.Module):
2     def __init__(self, dim, num_heads=8, qkv_bias=False, qk_scale=None, attn_drop=0., proj_drop=0.):
3         super().__init__()
4         self.num_heads = num_heads
5         head_dim = dim // num_heads
6         # NOTE scale factor was wrong in my original version, can set manually to be compat with prev
7         self.scale = qk_scale or head_dim ** -0.5
8
9         self.qkv = nn.Linear(dim, dim * 3, bias=qkv_bias)
10        self.attn_drop = nn.Dropout(attn_drop)
11        self.proj = nn.Linear(dim, dim)
12        self.proj_drop = nn.Dropout(proj_drop)
13
14    def forward(self, x):
15        B, N, C = x.shape
16        qkv = self.qkv(x).reshape(B, N, 3, self.num_heads, C // self.num_heads).permute(2, 0, 3, 1, 4)
17        q, k, v = qkv[0], qkv[1], qkv[2] # make torchscript happy (cannot use tensor as tuple)
18
19        attn = (q @ k.transpose(-2, -1)) * self.scale
20        attn = attn.softmax(dim=-1)
21        attn = self.attn_drop(attn)
22
23        x = (attn @ v).transpose(1, 2).reshape(B, N, C)
24        x = self.proj(x)
25        x = self.proj_drop(x)

```

```

26
27     # x: (B, N, C)
28     return x

```

6. 包含Attention和Add & Norm的Block实现：

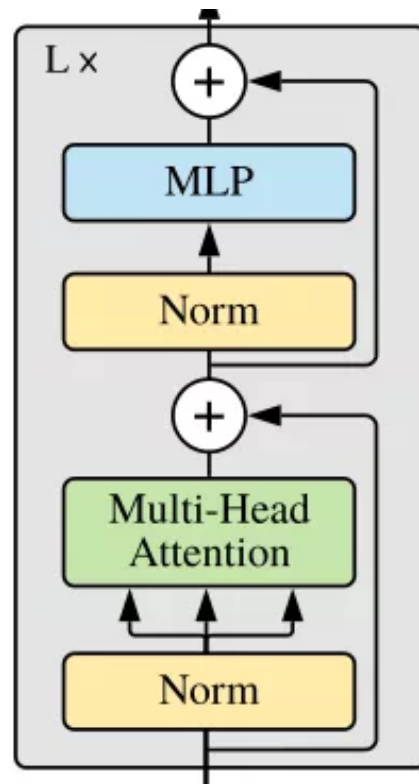


图1: Block类对应结构

不同之处是：

先进行Norm，再Attention；先进行Norm，再通过FFN (MLP)。

```

1 class Block(nn.Module):
2     def __init__(self, dim, num_heads, mlp_ratio=4., qkv_bias=False, qk_scale=None, drop=0., attn_drop
3         drop_path=0., act_layer=nn.GELU, norm_layer=nn.LayerNorm):
4         super().__init__()
5         self.norm1 = norm_layer(dim)
6         self.attn = Attention(
7             dim, num_heads=num_heads, qkv_bias=qkv_bias, qk_scale=qk_scale, attn_drop=attn_drop, proj_
8             # NOTE: drop path for stochastic depth, we shall see if this is better than dropout here
9             self.drop_path = DropPath(drop_path) if drop_path > 0. else nn.Identity()
10            self.norm2 = norm_layer(dim)
11            mlp_hidden_dim = int(dim * mlp_ratio)
12            self.mlp = Mlp(in_features=dim, hidden_features=mlp_hidden_dim, act_layer=act_layer, drop=drop
13
14    def forward(self, x):
15        x = x + self.drop_path(self.attn(self.norm1(x)))
16        x = x + self.drop_path(self.mlp(self.norm2(x)))
17        return x

```

7. 接下来要把图片转换成Patch，一种做法是直接把Image转化成Patch，另一种做法是把Backbone输出的特征转化成Patch。

1) 直接把Image转化成Patch:

输入的 x 的维度是: (B, C, H, W)

输出的PatchEmbedding的维度是: (B, 14*14, 768), 768表示embed_dim, 14*14表示一共有196个Patches。

```

1 class PatchEmbed(nn.Module):
2     """ Image to Patch Embedding
3     """
4     def __init__(self, img_size=224, patch_size=16, in_chans=3, embed_dim=768):
5         super().__init__()
6         img_size = to_2tuple(img_size)
7         patch_size = to_2tuple(patch_size)
8         num_patches = (img_size[1] // patch_size[1]) * (img_size[0] // patch_size[0])
9         self.img_size = img_size
10        self.patch_size = patch_size
11        self.num_patches = num_patches
12
13        self.proj = nn.Conv2d(in_chans, embed_dim, kernel_size=patch_size, stride=patch_size)
14
15    def forward(self, x):
16        B, C, H, W = x.shape
17        # FIXME look at relaxing size constraints
18        assert H == self.img_size[0] and W == self.img_size[1], \
19            f"Input image size ({H}*{W}) doesn't match model ({self.img_size[0]}*{self.img_size[1]})."
20        x = self.proj(x).flatten(2).transpose(1, 2)
21
22        # x: (B, 14*14, 768)
23        return x

```

2) 把Backbone输出的特征转化成Patch:

输入的 x 的维度是: (B, C, H, W)

得到Backbone输出的维度是: (B, feature_size, feature_size, feature_dim)

输出的PatchEmbedding的维度是: (B, feature_size, feature_size, embed_dim), 一共有feature_size * feature_size个Patches。

```

1 class HybridEmbed(nn.Module):
2     """ CNN Feature Map Embedding
3     Extract feature map from CNN, flatten, project to embedding dim.
4     """
5     def __init__(self, backbone, img_size=224, feature_size=None, in_chans=3, embed_dim=768):
6         super().__init__()
7         assert isinstance(backbone, nn.Module)
8         img_size = to_2tuple(img_size)
9         self.img_size = img_size
10        self.backbone = backbone
11        if feature_size is None:
12            with torch.no_grad():
13                # FIXME this is hacky, but most reliable way of determining the exact dim of the output

```

```

14         # map for all networks, the feature metadata has reliable channel and stride info, but
15         # stride to calc feature dim requires info about padding of each stage that isn't capt
16         training = backbone.training
17         if training:
18             backbone.eval()
19         o = self.backbone(torch.zeros(1, in_chans, img_size[0], img_size[1]))
20         if isinstance(o, (list, tuple)):
21             o = o[-1] # last feature if backbone outputs list/tuple of features
22         feature_size = o.shape[-2:]
23         feature_dim = o.shape[1]
24         backbone.train(training)
25     else:
26         feature_size = to_2tuple(feature_size)
27         if hasattr(self.backbone, 'feature_info'):
28             feature_dim = self.backbone.feature_info.channels()[-1]
29         else:
30             feature_dim = self.backbone.num_features
31         self.num_patches = feature_size[0] * feature_size[1]
32         self.proj = nn.Conv2d(feature_dim, embed_dim, 1)
33
34     def forward(self, x):
35         x = self.backbone(x)
36         if isinstance(x, (list, tuple)):
37             x = x[-1] # last feature if backbone outputs list/tuple of features
38         x = self.proj(x).flatten(2).transpose(1, 2)
39         return x

```

8. 以上是ViT所需的所有模块的定义，下面是VisionTransformer 这个类的实现：

8.1 使用这个类时需要传入的变量，其含义已经在本小节一开始介绍。

```

1 class VisionTransformer(nn.Module):
2     """ Vision Transformer
3
4     A PyTorch impl of : `An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale` -
5         https://arxiv.org/abs/2010.11929
6     """
7     def __init__(self, img_size=224, patch_size=16, in_chans=3, num_classes=1000, embed_dim=768, depth=
8         num_heads=12, mlp_ratio=4., qkv_bias=True, qk_scale=None, representation_size=None,
9         drop_rate=0., attn_drop_rate=0., drop_path_rate=0., hybrid_backbone=None, norm_layer=None)

```

8.2 得到分块后的Patch的数量：

```

1 super().__init__()
2 self.num_classes = num_classes
3 self.num_features = self.embed_dim = embed_dim # num_features for consistency with other models
4 norm_layer = norm_layer or partial(nn.LayerNorm, eps=1e-6)
5
6 if hybrid_backbone is not None:
7     self.patch_embed = HybridEmbed(
8         hybrid_backbone, img_size=img_size, in_chans=in_chans, embed_dim=embed_dim)
9 else:

```

```

10     self.patch_embed = PatchEmbed(
11         img_size=img_size, patch_size=patch_size, in_chans=in_chans, embed_dim=embed_dim)
12     num_patches = self.patch_embed.num_patches

```

8.3 class token:

一开始定义成(1, 1, 768)，之后再变成(B, 1, 768)。

```

1 self.cls_token = nn.Parameter(torch.zeros(1, 1, embed_dim))

```

8.4 定义位置编码:

```

1 self.pos_embed = nn.Parameter(torch.zeros(1, num_patches + 1, embed_dim))

```

8.5 把12个Block连接起来:

```

1 self.pos_drop = nn.Dropout(p=drop_rate)
2
3 dpr = [x.item() for x in torch.linspace(0, drop_path_rate, depth)] # stochastic depth decay rule
4 self.blocks = nn.ModuleList([
5     Block(
6         dim=embed_dim, num_heads=num_heads, mlp_ratio=mlp_ratio, qkv_bias=qkv_bias, qk_scale=qk_scale,
7         drop=drop_rate, attn_drop=attn_drop_rate, drop_path=dpr[i], norm_layer=norm_layer)
8     for i in range(depth)])
9 self.norm = norm_layer(embed_dim)

```

8.6 表示层和分类头:

表示层输出维度是representation_size，分类头输出维度是num_classes。

```

1 # Representation layer
2 if representation_size:
3     self.num_features = representation_size
4     self.pre_logits = nn.Sequential(OrderedDict([
5         ('fc', nn.Linear(embed_dim, representation_size)),
6         ('act', nn.Tanh())
7     ]))
8 else:
9     self.pre_logits = nn.Identity()
10
11 # Classifier head
12 self.head = nn.Linear(self.num_features, num_classes) if num_classes > 0 else nn.Identity()

```

8.7 初始化各个模块:

函数trunc_normal_(tensor, mean=0., std=1., a=-2., b=2.)的目的是用截断的正态分布绘制的值填充输入张量，我们只需要输入均值mean，标准差std，下界a，上界b即可。

`self.apply(self._init_weights)`表示对各个模块的权重进行初始化。`apply`函数的代码是：

```
1     for module in self.children():
2         module.apply(fn)
3     fn(self)
4     return self
```

递归地将`fn`应用于每个子模块，相当于在递归调用`fn`，即`_init_weights`这个函数。
也就是把模型的所有子模块的`nn.Linear`和`nn.LayerNorm`层都初始化掉。

```
1 trunc_normal_(self.pos_embed, std=.02)
2 trunc_normal_(self.cls_token, std=.02)
3 self.apply(self._init_weights)
4
5 def _init_weights(self, m):
6     if isinstance(m, nn.Linear):
7         trunc_normal_(m.weight, std=.02)
8         if isinstance(m, nn.Linear) and m.bias is not None:
9             nn.init.constant_(m.bias, 0)
10    elif isinstance(m, nn.LayerNorm):
11        nn.init.constant_(m.bias, 0)
12        nn.init.constant_(m.weight, 1.0)
```

8.8 最后就是整个ViT模型的forward实现：

```
1 def forward_features(self, x):
2     B = x.shape[0]
3     x = self.patch_embed(x)
4
5     cls_tokens = self.cls_token.expand(B, -1, -1) # stole cls_tokens impl from Phil Wang, thanks
6     x = torch.cat((cls_tokens, x), dim=1)
7     x = x + self.pos_embed
8     x = self.pos_drop(x)
9
10    for blk in self.blocks:
11        x = blk(x)
12
13    x = self.norm(x)[:, 0]
14    x = self.pre_logits(x)
15    return x
16
17 def forward(self, x):
18     x = self.forward_features(x)
19     x = self.head(x)
20     return x
```

9. 下面是Training data-efficient image transformers & distillation through attention这篇论文的DeiT这个类的实现：

整体结构与ViT相似，继承了上面的VisionTransformer类。

```
1 class DistilledVisionTransformer(VisionTransformer):
```

再额外定义以下3个变量：

- **distillation token: dist_token**
- **新的位置编码: pos_embed**
- **蒸馏分类头: head_dist**

DeiT相关介绍可以参考：Vision Transformer 超详细解读 (原理分析+代码解读) (三)。

```
1 self.dist_token = nn.Parameter(torch.zeros(1, 1, self.embed_dim))
2 num_patches = self.patch_embed.num_patches
3 self.pos_embed = nn.Parameter(torch.zeros(1, num_patches + 2, self.embed_dim))
4 self.head_dist = nn.Linear(self.embed_dim, self.num_classes) if self.num_classes > 0 else nn.Identity()
```

初始化新定义的变量：

```
1 trunc_normal_(self.dist_token, std=.02)
2 trunc_normal_(self.pos_embed, std=.02)
3 self.head_dist.apply(self._init_weights)
```

前向函数：

```
1 def forward_features(self, x):
2     B = x.shape[0]
3     x = self.patch_embed(x)
4
5     cls_tokens = self.cls_token.expand(B, -1, -1) # stole cls_tokens impl from Phil Wang, thanks
6     dist_token = self.dist_token.expand(B, -1, -1)
7     x = torch.cat((cls_tokens, dist_token, x), dim=1)
8
9     x = x + self.pos_embed
10    x = self.pos_drop(x)
11
12    for blk in self.blocks:
13        x = blk(x)
14
15    x = self.norm(x)
16    return x[:, 0], x[:, 1]
17
18 def forward(self, x):
19    x, x_dist = self.forward_features(x)
20    x = self.head(x)
21    x_dist = self.head_dist(x_dist)
22    if self.training:
23        return x, x_dist
24    else:
25        # during inference, return the average of both classifier predictions
```

```
26         return (x + x_dist) / 2
```

10. 对位置编码进行插值：

posemb代表未插值的位置编码权值，posemb_tok为位置编码的token部分，posemb_grid为位置编码的插值部分。

首先把要插值部分posemb_grid给reshape成(1, gs_old, gs_old, -1)的形式，再插值成(1, gs_new, gs_new, -1)的形式，最后与token部分在第1维度拼接在一起，得到插值后的位置编码posemb。

```
1 def resize_pos_embed(posemb, posemb_new):
2     # Rescale the grid of position embeddings when loading from state_dict. Adapted from
3     # https://github.com/google-research/vision_transformer/blob/00883dd691c63a6830751563748663526e811
4     _logger.info('Resized position embedding: %s to %s', posemb.shape, posemb_new.shape)
5     ntok_new = posemb_new.shape[1]
6     if True:
7         posemb_tok, posemb_grid = posemb[:, :1], posemb[0, 1:]
8         ntok_new -= 1
9     else:
10        posemb_tok, posemb_grid = posemb[:, :0], posemb[0]
11    gs_old = int(math.sqrt(len(posemb_grid)))
12    gs_new = int(math.sqrt(ntok_new))
13    _logger.info('Position embedding grid-size from %s to %s', gs_old, gs_new)
14    posemb_grid = posemb_grid.reshape(1, gs_old, gs_old, -1).permute(0, 3, 1, 2)
15    posemb_grid = F.interpolate(posemb_grid, size=(gs_new, gs_new), mode='bilinear')
16    posemb_grid = posemb_grid.permute(0, 2, 3, 1).reshape(1, gs_new * gs_new, -1)
17    posemb = torch.cat([posemb_tok, posemb_grid], dim=1)
18    return posemb
```

11. _create_vision_transformer函数用于创建vision transformer：

checkpoint_filter_fn的作用是加载预训练权重。

```
1 def checkpoint_filter_fn(state_dict, model):
2     """ convert patch embedding weight from manual patchify + linear proj to conv"""
3     out_dict = {}
4     if 'model' in state_dict:
5         # For deit models
6         state_dict = state_dict['model']
7     for k, v in state_dict.items():
8         if 'patch_embed.proj.weight' in k and len(v.shape) < 4:
9             # For old models that I trained prior to conv based patchification
10            0, I, H, W = model.patch_embed.proj.weight.shape
11            v = v.reshape(0, -1, H, W)
12        elif k == 'pos_embed' and v.shape != model.pos_embed.shape:
13            # To resize pos embedding when using model at different size from pretrained weights
14            v = resize_pos_embed(v, model.pos_embed)
15        out_dict[k] = v
16    return out_dict
17
18
19 def _create_vision_transformer(variant, pretrained=False, distilled=False, **kwargs):
```

```

20     default_cfg = default_cfgs[variant]
21     default_num_classes = default_cfg['num_classes']
22     default_img_size = default_cfg['input_size'][-1]
23
24     num_classes = kwargs.pop('num_classes', default_num_classes)
25     img_size = kwargs.pop('img_size', default_img_size)
26     repr_size = kwargs.pop('representation_size', None)
27     if repr_size is not None and num_classes != default_num_classes:
28         # Remove representation layer if fine-tuning. This may not always be the desired action,
29         # but I feel better than doing nothing by default for fine-tuning. Perhaps a better interface?
30         _logger.warning("Removing representation layer for fine-tuning.")
31         repr_size = None
32
33     model_cls = DistilledVisionTransformer if distilled else VisionTransformer
34     model = model_cls(img_size=img_size, num_classes=num_classes, representation_size=repr_size, **kwargs)
35     model.default_cfg = default_cfg
36
37     if pretrained:
38         load_pretrained(
39             model, num_classes=num_classes, in_chans=kwargs.get('in_chans', 3),
40             filter_fn=partial(checkpoint_filter_fn, model=model))
41     return model

```

12. 定义和注册vision transformer模型：

@ 指装饰器。

@register_model代表注册器，注册这个新定义的模型。

model_kwargs是一个存有模型所有超参数的字典。

最后使用上面定义的_create_vision_transformer函数创建模型。

```

1 @register_model
2 def vit_base_patch16_224(pretrained=False, **kwargs):
3     """ ViT-Base (ViT-B/16) from original paper (https://arxiv.org/abs/2010.11929).
4     ImageNet-1k weights fine-tuned from in21k @ 224x224, source https://github.com/google-research/vision_transformer
5     """
6     model_kwargs = dict(patch_size=16, embed_dim=768, depth=12, num_heads=12, **kwargs)
7     model = _create_vision_transformer('vit_base_patch16_224', pretrained=pretrained, **model_kwargs)
8     return model

```

一共可以选择的模型包括：

ViT系列：

```

vit_small_patch16_224
vit_base_patch16_224
vit_base_patch32_224
vit_base_patch16_384
vit_base_patch32_384
vit_large_patch16_224
vit_large_patch32_224
vit_large_patch16_384

```

```
vit_large_patch32_384
vit_base_patch16_224_in21k
vit_base_patch32_224_in21k
vit_large_patch16_224_in21k
vit_large_patch32_224_in21k
vit_huge_patch14_224_in21k
vit_base_resnet50_224_in21k
vit_base_resnet50_384
vit_small_resnet26d_224
vit_small_resnet50d_s3_224
vit_base_resnet26d_224
vit_base_resnet50d_224
```

DeiT系列：

```
vit_deit_tiny_patch16_224
vit_deit_small_patch16_224
vit_deit_base_patch16_224
vit_deit_base_patch16_384
vit_deit_tiny_distilled_patch16_224
vit_deit_small_distilled_patch16_224
vit_deit_base_distilled_patch16_224
vit_deit_base_distilled_patch16_384
```

以上就是对timm库 vision_transformer.py代码的分析。

如何使用timm库以及 vision_transformer.py代码搭建自己的模型？

在搭建我们自己的视觉Transformer模型时，我们可以按照下面的步骤操作：首先

- 继承timm库的**VisionTransformer**这个类。
- 添加上自己模型独有的一些变量。
- 重写**forward**函数。
- 通过timm库的注册器注册新模型。

我们以ViT模型的改进版DeiT为例：

首先，DeiT的所有模型列表如下：

```
1 __all__ = [
2     'deit_tiny_patch16_224', 'deit_small_patch16_224', 'deit_base_patch16_224',
3     'deit_tiny_distilled_patch16_224', 'deit_small_distilled_patch16_224',
4     'deit_base_distilled_patch16_224', 'deit_base_patch16_384',
5     'deit_base_distilled_patch16_384',
6 ]
```

导入VisionTransformer这个类，注册器register_model，以及初始化函数trunc_normal_：

```
1 from timm.models.vision_transformer import VisionTransformer, _cfg
2 from timm.models.registry import register_model
3 from timm.models.layers import trunc_normal_
```


DeiT的class名称是DistilledVisionTransformer，它直接继承了VisionTransformer这个类：

```
1 class DistilledVisionTransformer(VisionTransformer):
```

添加上自己模型独有的一些变量：

```
1 def __init__(self, *args, **kwargs):
2     super().__init__(*args, **kwargs)
3     self.dist_token = nn.Parameter(torch.zeros(1, 1, self.embed_dim))
4     num_patches = self.patch_embed.num_patches
5     # 位置编码不是ViT中的(b, N, 256)，而变成了(b, N+2, 256)，原因是还有class token和distillation token.
6     self.pos_embed = nn.Parameter(torch.zeros(1, num_patches + 2, self.embed_dim))
7     self.head_dist = nn.Linear(self.embed_dim, self.num_classes) if self.num_classes > 0 else nn.Identity()
8
9     trunc_normal_(self.dist_token, std=.02)
10    trunc_normal_(self.pos_embed, std=.02)
11    self.head_dist.apply(self._init_weights)
```

重写forward函数：

```
1 def forward_features(self, x):
2     # taken from https://github.com/rwightman/pytorch-image-models/blob/master/timm/models/vision_transformer.py
3     # with slight modifications to add the dist_token
4     B = x.shape[0]
5
6     x = self.patch_embed(x)
7
8     cls_tokens = self.cls_token.expand(B, -1, -1) # stole cls_tokens impl from Phil Wang, thanks
9     dist_token = self.dist_token.expand(B, -1, -1)
10
11    x = torch.cat((cls_tokens, dist_token, x), dim=1)
12
13    x = x + self.pos_embed
14    x = self.pos_drop(x)
15
16    for blk in self.blocks:
17        x = blk(x)
18
19    x = self.norm(x)
20
21    return x[:, 0], x[:, 1]
22
23 def forward(self, x):
24    x, x_dist = self.forward_features(x)
25    x = self.head(x)
26    x_dist = self.head_dist(x_dist)
27    if self.training:
28        return x, x_dist
29    else:
30        # during inference, return the average of both classifier predictions
31        return (x + x_dist) / 2
```

通过timm库的注册器注册新模型：

```
1 @register_model
2 def deit_base_patch16_224(pretrained=False, **kwargs):
3     model = VisionTransformer(
4         patch_size=16, embed_dim=768, depth=12, num_heads=12, mlp_ratio=4, qkv_bias=True,
5         norm_layer=partial(nn.LayerNorm, eps=1e-6), **kwargs)
6     model.default_cfg = _cfg()
7     if pretrained:
8         checkpoint = torch.hub.load_state_dict_from_url(
9             url="https://dl.fbaipublicfiles.com/deit/deit_base_patch16_224-b5f2ef4d.pth",
10            map_location="cpu", check_hash=True
11        )
12        model.load_state_dict(checkpoint["model"])
13    return model
```

推荐阅读

一文看懂 9 种Transformer结构！



2021-02-18

来自Transformer的降维打击：ReID各项任务全面领先，阿里&浙大提出TransReID



2021-02-09

网络架构设计：CNN based和Transformer based



2021-02-08



#原创作者激励计划#

极市平台深耕CV开发者领域近5年，拥有一大批优质CV开发者受众，覆盖微信、知乎、B站、微博等多个渠道。通过极市平台，您的文章的观点和看法能分享至更多CV开发者，既能体现文章的价值，又能让文章在视觉圈内得到更大程度上的推广。

对于优质内容开发者，极市可推荐至国内优秀出版社作出书，同时为开发者引荐行业大牛，组织个人分享交流会，推荐名企就业机会，打造个人品牌 IP。

投稿须知：

- 1.作者保证投稿作品为自己的原创作品。
- 2.极市平台尊重原作者署名权，并支付相应稿费。文章发布后，版权仍属于原作者。
- 3.原作者可以将文章发在其他平台的个人账号，但需要在文章顶部标明首发于极市平台

投稿方式：

添加小编微信Fengcall（微信号：fengcall19），备注：姓名-投稿



△长按添加极市平台小编

△点击卡片关注极市平台，获取**最新CV干货**

觉得有用麻烦给个在看啦~ 

阅读原文

喜欢此内容的人还喜欢

15个目标检测开源数据集汇总
极市平台