

```
1 import torch
2 import torch.nn as nn
3 import torchvision.models as models
4 from torch.nn.utils.rnn import pack_padded_sequence
5 import torch.nn.functional as F
6 import numpy as np
7 import pdb
8 from .densenet import *
9 from .resnet import *
10 from .vgg import *
11 import sys
12 thismodule = sys.modules[__name__]
13
14
15 dim_dict = {
16     'resnet101': [512, 1024, 2048],
17     'resnet152': [512, 1024, 2048],
18     'resnet50': [512, 1024, 2048],
19     'resnet34': [128, 256, 512],
20     'resnet18': [128, 256, 512],
21     'densenet121': [256, 512, 1024],
22     'densenet161': [384, 1056, 2208],
23     'densenet169': [64, 128, 256, 640, 1664],
24     # 'densenet169': [256, 640, 1664],
25     'densenet201': [256, 896, 1920],
26     'vgg': [256, 512, 512]
27 }
28
29
30 def get_upsampling_weight(in_channels, out_channels, kernel_size):
31     """Make a 2D bilinear kernel suitable for upsampling"""
32     factor = (kernel_size + 1) // 2
33     if kernel_size % 2 == 1:
34         center = factor - 1
35     else:
36         center = factor - 0.5
37     og = np.ogrid[:kernel_size, :kernel_size]
38     filt = (1 - abs(og[0] - center) / factor) * \
39           (1 - abs(og[1] - center) / factor)
40     weight = np.zeros((in_channels, out_channels, kernel_size,
41                        kernel_size),
42                       dtype=np.float64)
43     weight[range(in_channels), range(out_channels), :, :] = filt
44     return torch.from_numpy(weight).float()
45
46 class ParamPool(nn.Module):
47     def __init__(self, input_c):
48         super(ParamPool, self).__init__()
49         self.conv = nn.Conv2d(input_c, 1, kernel_size=1, bias=False)
50
51     def forward(self, x):
52         bsize, c, ssize, _ = x.shape
53         w = self.conv(x)
54         w = F.softmax(w.view(bsize, 1, -1), 2)
55         w = w.view(bsize, 1, ssize, ssize)
56         x = (x*w).sum(3).sum(2)
57         return x
58
59
```

```

60 def proc_densenet(model):
61     def hook(module, input, output):
62         model.feats[output.device.index] += [output]
63     model.features.transition3[-2].register_forward_hook(hook)
64     model.features.transition2[-2].register_forward_hook(hook)
65     # dilation
66     def remove_sequential(all_layers, network):
67         for layer in network.children():
68             if isinstance(layer, nn.Sequential): # if sequential layer,
        apply recursively to layers in sequential layer
69                 remove_sequential(all_layers, layer)
70             if list(layer.children()) == []: # if leaf node, add it to
        list
71                 all_layers.append(layer)
72     model.features.transition3[-1].kernel_size = 1
73     model.features.transition3[-1].stride = 1
74     all_layers = []
75     remove_sequential(all_layers, model.features.denseblock4)
76     for m in all_layers:
77         if isinstance(m, nn.Conv2d) and m.kernel_size==(3, 3):
78             m.dilation = (2, 2)
79             m.padding = (2, 2)
80     return model
81
82
83 procs = {
84     'densenet169': proc_densenet,
85     'densenet201': proc_densenet,
86 }
87
88
89 class EncoderCNN(nn.Module):
90     def __init__(self, embed_size, patt_size=512, base='densenet169'):
91         """Load the pretrained ResNet-152 and replace top fc layer."""
92         super(EncoderCNN, self).__init__()
93         if 'vgg' in base:
94             dims = dim_dict['vgg'][::-1]
95         else:
96             dims = dim_dict[base][::-1]
97         self.preds = nn.ModuleList([nn.Conv2d(d, 1, kernel_size=1) for d
        in dims])
98         self.upscales = nn.ModuleList([
99             nn.ConvTranspose2d(1, 1, 4, 2, 1),
100             nn.ConvTranspose2d(1, 1, 4, 2, 1),
101             nn.ConvTranspose2d(1, 1, 16, 8, 4),
102         ])
103         self.linear = nn.Linear(patt_size, embed_size)
104         self.reduce = nn.Conv2d(dims[0], patt_size, kernel_size=3,
        padding=1)
105         self.param_pool = ParamPool(patt_size)
106         for m in self.modules():
107             if isinstance(m, nn.Conv2d) or isinstance(m, nn.Linear):
108                 m.weight.data.normal_(0.0, 0.02)
109                 if m.bias is not None:
110                     m.bias.data.fill_(0)
111             if isinstance(m, nn.ConvTranspose2d):
112                 assert m.kernel_size[0] == m.kernel_size[1]
113                 initial_weight = get_upsampling_weight(
114                     m.in_channels, m.out_channels, m.kernel_size[0])
115                 m.weight.data.copy_(initial_weight)

```

```

116
117     self.feature = getattr(thismodule, base)(pretrained=True)
118     self.feature.feats = {}
119     self.feature = procs[base](self.feature)
120     for m in self.feature.modules():
121         if isinstance(m, nn.BatchNorm2d):
122             m.requires_grad=False
123
124     def forward(self, x):
125         """Extract feature vectors from input images."""
126         # pdb.set_trace()
127         # with torch.no_grad():
128             features = self.resnet(images)
129         # features = features.reshape(features.size(0), -1)
130         # features = self.bn(self.linear(features))
131         # return features
132         self.feature.feats[x.device.index] = []
133         x = self.feature(x)
134         feats = self.feature.feats[x.device.index]
135         feats += [x]
136         feats = feats[:-1]
137         msk = self.preds[0](feats[0])
138         big_msk = msk
139         # big_msk = self.upscales[0](msk)
140         # big_msk = F.upsample(msk, scale_factor=16)
141         msk = F.sigmoid(msk)
142         msk_feat = self.reduce(feats[0])*msk
143         # msk_feat = F.avg_pool2d(msk_feat, 16).squeeze(3).squeeze(2)
144         msk_feat = self.param_pool(msk_feat)
145         msk_feat = self.linear(msk_feat)
146         big_msk = F.upsample_bilinear(big_msk, scale_factor=16)
147         if self.training:
148             return big_msk, msk, msk_feat
149         else:
150             return big_msk
151
152
153 class DecoderRNN(nn.Module):
154     def __init__(self, embed_size, hidden_size, vocab_size, num_layers,
155 max_seq_length=20):
156         """Set the hyper-parameters and build the layers."""
157         super(DecoderRNN, self).__init__()
158         self.embed = nn.Embedding(vocab_size, embed_size)
159         self.lstm = nn.LSTM(embed_size, hidden_size, num_layers,
160 batch_first=True)
161         self.linear = nn.Linear(hidden_size, vocab_size)
162         self.max_seq_length = max_seq_length
163
164     def forward(self, features, captions, lengths):
165         """Decode image feature vectors and generates captions."""
166         embeddings = self.embed(captions)
167         embeddings = torch.cat((features.unsqueeze(1), embeddings), 1)
168         packed = pack_padded_sequence(embeddings, lengths,
169 batch_first=True)
170         hiddens, _ = self.lstm(packed)
171         outputs = self.linear(hiddens[0])
172         return outputs
173
174     def sample(self, features, states=None):
175         """Generate captions for given image features using greedy

```

```
search. """
173     sampled_ids = []
174     inputs = features.unsqueeze(1)
175     for i in range(self.max_seq_length):
176         hiddens, states = self.lstm(inputs, states) #
177         hiddens: (batch_size, 1, hidden_size)
178         outputs = self.linear(hiddens.squeeze(1)) #
179         outputs: (batch_size, vocab_size)
180         _, predicted = outputs.max(1) #
181         predicted: (batch_size)
182         sampled_ids.append(predicted)
183         inputs = self.embed(predicted) #
184         inputs: (batch_size, embed_size)
185         inputs = inputs.unsqueeze(1) #
186         inputs: (batch_size, 1, embed_size)
187         sampled_ids = torch.stack(sampled_ids, 1) #
188         sampled_ids: (batch_size, max_seq_length)
189         return sampled_ids
190
191
192 class EncDec(nn.Module):
193     def __init__(self, vocab_size, embed_size=256, hidden_size=512,
194                 num_layers=1, max_seq_length=20):
195         super(EncDec, self).__init__()
196         self.encoder = EncoderCNN(embed_size)
197         self.decoder = DecoderRNN(embed_size, hidden_size, vocab_size,
198                                 num_layers, max_seq_length)
199
200     def forward(self, images, captions=None, lengths=None):
201         if self.training:
202             if captions is not None:
203                 big_msk, msk, msk_feat = self.encoder(images)
204                 outputs = self.decoder(msk_feat, captions, lengths)
205                 return big_msk, msk, outputs
206             else:
207                 big_msk, msk, msk_feat = self.encoder(images)
208                 return big_msk, msk
209         else:
210             return self.encoder(images)
```