Homework 4

Instructions

- This homework focuses on understanding and applying CoCoOp for CLIP prompt tuning. It consists of **four questions** designed to assess both theoretical understanding and practical application.
- · Please organize your answers and results for the questions below and submit this jupyter notebook as a .pdf file.
- Deadline: 11/26 (Sat) 23:59

> Preparation

- Run the code below before proceeding with the homework.
- · If an error occurs, click 'Run Session Again' and then restart the runtime from the beginning.



Q1. Understanding and implementing CoCoOp

- · We have learned how to define CoOp in Lab Session 4.
- The main difference between CoOp and CoCoOp is meta network to extract image tokens that is added to the text prompt.
- Based on the CoOp code given in Lab Session 4, fill-in-the-blank exercise (4 blanks!!) to test your understanding of critical parts of the CoCoOp.

```
1 import torch.nn as nn
3 class CoCoOpPromptLearner(nn.Module):
 4
      def __init__(self, cfg, classnames, clip_model):
5
          super().__init__()
 6
          n_cls = len(classnames)
 7
          n_ctx = cfg.TRAINER.COCOOP.N_CTX
8
          ctx_init = cfg.TRAINER.COCOOP.CTX_INIT
9
          dtype = clip_model.dtype
          ctx_dim = clip_model.ln_final.weight.shape[0]
10
11
          vis_dim = clip_model.visual.output_dim
12
          clip_imsize = clip_model.visual.input_resolution
          cfq imsize = cfq.INPUT.SIZE[0]
13
14
          assert cfg_imsize == clip_imsize, f"cfg_imsize ({cfg_imsize}) must equal to clip_imsize ({clip_imsize})"
15
16
          if ctx_init:
17
              # use given words to initialize context vectors
18
              ctx_init = ctx_init.replace("
              n_ctx = len(ctx_init.split(" "))
19
              prompt = clip.tokenize(ctx_init)
20
21
              with torch.no_grad():
22
                  embedding = clip_model.token_embedding(prompt).type(dtype)
23
              ctx_vectors = embedding[0, 1: 1 + n_ctx, :]
24
              prompt_prefix = ctx_init
25
          else:
26
              # random initialization
27
              ctx_vectors = torch.empty(n_ctx, ctx_dim, dtype=dtype)
28
              nn.init.normal_(ctx_vectors, std=0.02)
              prompt_prefix = " ".join(["X"] * n_ctx)
29
30
          print(f'Initial context: "{prompt_prefix}"')
31
32
          print(f"Number of context words (tokens): {n_ctx}")
33
          self.ctx = nn.Parameter(ctx_vectors) # Wrap the initialized prompts above as parameters to make them trainable.
34
35
          ### Tokenize ###
36
          classnames = [name.replace("_", " ") for name in classnames] # 예) "Forest"
37
38
          name_lens = [len(_tokenizer.encode(name)) for name in classnames]
          prompts = [prompt_prefix + " " + name + "." for name in classnames] # 예) "A photo of Forest."
39
40
          tokenized_prompts = torch.cat([clip.tokenize(p) for p in prompts]) # 예) [49406, 320, 1125, 539...]
41
42
43
44
45
          46
          ###### Q1. Fill in the blank ######
47
          ######## Define Meta Net ########
48
          self.meta_net = nn.Sequential(OrderedDict([
              ("linear1", nn.Linear(vis_dim, vis_dim // 16)),
49
```

```
50
               ("relu", nn.ReLU(inplace=Irue)),
 51
               ("linear2", nn.Linear(vis_dim // 16, ctx_dim))
           1))
 52
 53
           54
           ## Hint: meta network is composed to linear layer, relu activation, and linear layer.
 55
 56
 57
 58
           if cfg.TRAINER.COCOOP.PREC == "fp16":
 59
               self.meta_net.half()
 60
 61
           with torch.no_grad():
               embedding = clip_model.token_embedding(tokenized_prompts).type(dtype)
 62
 63
 64
           # These token vectors will be saved when in save_model(),
 65
           # but they should be ignored in load model() as we want to use
 66
           # those computed using the current class names
           self.register_buffer("token_prefix", embedding[:, :1, :]) # SOS
self.register_buffer("token_suffix", embedding[:, 1 + n_ctx:, :]) # CLS, EOS
 67
 68
 69
           self.n_cls = n_cls
 70
           self.n ctx = n ctx
 71
           self.tokenized_prompts = tokenized_prompts # torch.Tensor
 72
           self.name_lens = name_lens
 73
 74
       def construct_prompts(self, ctx, prefix, suffix, label=None):
 75
           # dim0 is either batch_size (during training) or n_cls (during testing)
 76
           # ctx: context tokens, with shape of (dim0, n_ctx, ctx_dim)
 77
           # prefix: the sos token, with shape of (n_cls, 1, ctx_dim)
 78
           # suffix: remaining tokens, with shape of (n_cls, *, ctx_dim)
 79
 80
           if label is not None:
 81
               prefix = prefix[label]
               suffix = suffix[label]
 82
 83
 84
           prompts = torch.cat(
 85
               [
                   prefix, \# (dim0, 1, dim)
 86
 87
                   ctx, # (dim0, n_ctx, dim)
 88
                   suffix, # (dim0, *, dim)
 89
               ],
 90
               dim=1,
 91
           )
 92
 93
           return prompts
 94
 95
       def forward(self, im_features):
 96
           prefix = self.token_prefix
 97
           suffix = self.token_suffix
 98
           ctx = self.ctx # (n_ctx, ctx_dim)
 99
100
101
102
           103
           ######## Q2,3. Fill in the blank #######
           bias = self.meta_net(im_features) # (batch, ctx_dim)
104
105
           bias = bias.unsqueeze(1) # (batch, 1, ctx_dim)
106
           ctx = ctx.unsqueeze(0) # (1, n_ctx, ctx_dim)
           ctx_shifted = ctx + bias # (batch, n_ctx, ctx_dim)
107
           108
109
           110
111
112
           # Use instance-conditioned context tokens for all classes
113
           prompts = []
114
115
           for ctx_shifted_i in ctx_shifted:
               ctx_i = ctx_shifted_i.unsqueeze(0).expand(self.n_cls, -1, -1)
116
               pts\_i = self.construct\_prompts(ctx\_i, prefix, suffix) \quad \# \ (n\_cls, n\_tkn, \ ctx\_dim)
117
118
               prompts.append(pts_i)
119
           prompts = torch.stack(prompts)
120
121
           return prompts
 1 class CoCoOpCustomCLIP(nn.Module):
 2
       def __init__(self, cfg, classnames, clip_model):
 3
           super().__init__()
           self.prompt_learner = CoCoOpPromptLearner(cfg, classnames, clip_model)
  4
  5
           self.tokenized_prompts = self.prompt_learner.tokenized_prompts
  6
           self.image_encoder = clip_model.visual
  7
           self.text_encoder = TextEncoder(clip_model)
  8
           self.logit_scale = clip_model.logit_scale
           self.dtype = clip_model.dtype
```

```
10
      def forward(self, image, label=None):
11
12
         tokenized prompts = self.tokenized prompts
13
         logit_scale = self.logit_scale.exp()
14
15
         image_features = self.image_encoder(image.type(self.dtype))
16
         image_features = image_features / image_features.norm(dim=-1, keepdim=True)
17
18
         19
         ######## 04. Fill in the blank #######
20
21
         prompts = self.prompt_learner(image_features)
22
         23
         24
25
26
         logits = []
27
         for pts_i, imf_i in zip(prompts, image_features):
28
            text_features = self.text_encoder(pts_i, tokenized_prompts)
29
             text_features = text_features / text_features.norm(dim=-1, keepdim=True)
30
             l_i = logit_scale * imf_i @ text_features.t()
31
             logits.append(l_i)
32
         logits = torch.stack(logits)
33
34
         if self.prompt_learner.training:
35
             return F.cross_entropy(logits, label)
36
         return logits
37
```

∨ Q2. Training CoCoOp

In this task, you will train CoCoOp on the EuroSAT dataset. If your implementation of CoCoOp in Question 1 is correct, the following code should execute without errors. Please submit the execution file so we can evaluate whether your code runs without any issues.

```
1 # Train on the Base Classes Train split and evaluate accuracy on the Base Classes Test split.
2 args.trainer = "CoCoOp"
3 args.train_batch_size = 4
4 args.epoch = 100
5 args.output_dir = "outputs/cocoop"
7 args.subsample_classes = "base"
8 args.eval only = False
9 cocoop_base_acc = main(args)
→ Loading trainer: CoCoOp
    Loading dataset: EuroSAT
    Reading split from /content/ProMetaR/data/eurosat/split_zhou_EuroSAT.json
    Creating a 16-shot dataset
    Creating a 4-shot dataset
    Saving preprocessed few-shot data to /content/ProMetaR/data/eurosat/split_fewshot/shot_16-seed_1.pkl
    SUBSAMPLE BASE CLASSES!
    Building transform_train
    + random resized crop (size=(224, 224), scale=(0.08, 1.0))
    + random flip
    + to torch tensor of range [0, 1]
    + normalization (mean=[0.48145466, 0.4578275, 0.40821073], std=[0.26862954, 0.26130258, 0.27577711])
    Building transform_test
    + resize the smaller edge to 224
    + 224x224 center crop
    + to torch tensor of range [0, 1]
    + normalization (mean=[0.48145466, 0.4578275, 0.40821073], std=[0.26862954, 0.26130258, 0.27577711])
    Dataset
                EuroSAT
    # classes
    # train_x
                20
    # val
    # test
                4,200
    Loading CLIP (backbone: ViT-B/16)
    /usr/local/lib/python3.10/dist-packages/torch/utils/data/dataloader.py:617: UserWarning: This DataLoader will create 8 w
      warnings.warn(
    Building custom CLIP
    Initial context: "a photo of a"
    Number of context words (tokens): 4
    Turning off gradients in both the image and the text encoder
    Parameters to be updated: {'prompt_learner.ctx', 'prompt_learner.meta_net.linear2.weight', 'prompt_learner.meta_net.line
    Loading evaluator: Classification
    No checkpoint found, train from scratch
    /usr/local/lib/python3.10/dist-packages/torch/optim/lr_scheduler.py:62: UserWarning: The verbose parameter is deprecated
      warnings.warn(
    epoch [1\bar{1}]100] batch [20/20] time 0.138 (0.353) data 0.000 (0.032) loss 0.2744 (1.1881) lr 2.5000e-03 eta 0:11:38 epoch [2/100] batch [20/20] time 0.104 (0.129) data 0.000 (0.018) loss 0.8384 (0.8970) lr 2.4994e-03 eta 0:04:11
    epoch [3/100] batch [20/20] time 0.095 (0.131) data 0.000 (0.017) loss 0.6382 (0.7859) lr 2.4975e-03 eta 0:04:14
```

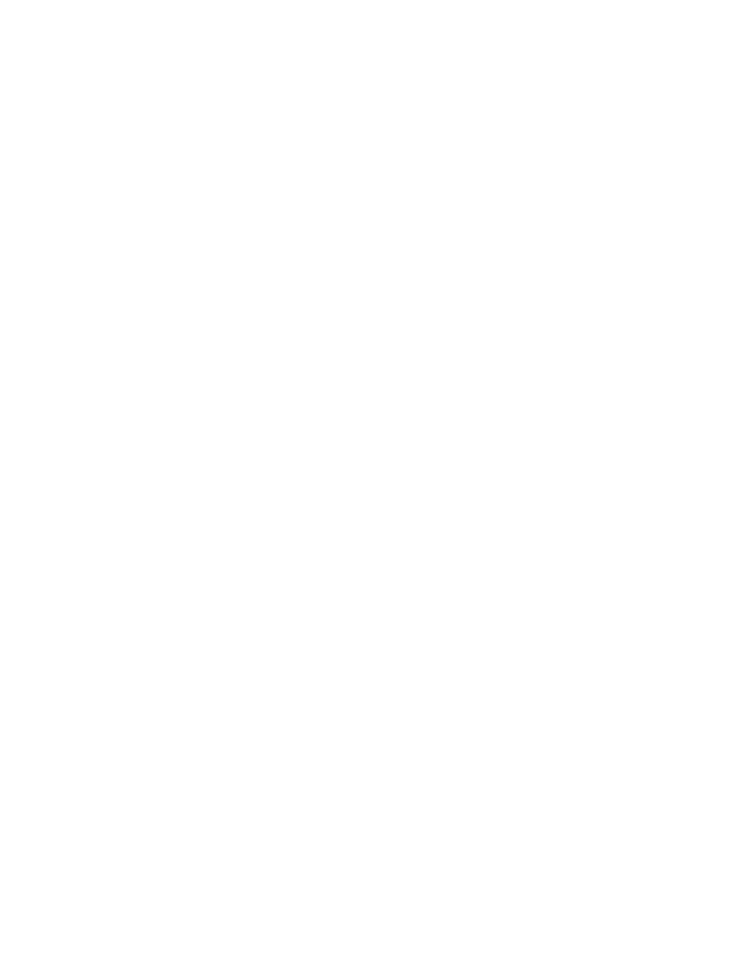
```
epoch [4/100] batch [20/20] time 0.095 (0.127) data 0.000 (0.016) loss 0.5044 (0.7151) lr 2.4945e-03 eta 0:04:03
    epoch [5/100]
                  batch [20/20] time 0.121 (0.142) data 0.000 (0.015) loss 0.5703 (0.6317) lr 2.4901e-03 eta 0:04:29
    epoch [6/100]
                  batch [20/20] time 0.154 (0.192) data 0.000 (0.034)
                                                                      loss 0.6060 (0.6009) lr 2.4846e-03 eta 0:06:01
                                                              (0.022)
    epoch [7/100]
                  batch [20/20] time 0.097
                                           (0.136)
                                                   data 0.000
                                                                       loss 0.3853 (0.6638) lr 2.4779e-03 eta 0:04:12
    epoch [8/100]
                  batch [20/20]
                               time 0.101 (0.136) data 0.000
                                                              (0.019)
                                                                      loss 1.4082 (0.6633) lr 2.4699e-03 eta 0:04:09
    epoch [9/100] batch [20/20] time 0.095 (0.131) data 0.000 (0.022) loss 0.1780 (0.4582) lr 2.4607e-03 eta 0:03:58
    epoch [10/100] batch [20/20] time 0.141 (0.152) data 0.000 (0.016) loss 1.2285 (0.5051) lr 2.4504e-03 eta 0:04:33
    epoch [11/100] batch [20/20] time 0.106 (0.199) data 0.000 (0.036) loss 0.2539 (0.5013) lr 2.4388e-03 eta 0:05:54
                         [20/20]
                                 time 0.094 (0.126)
                                                               (0.015)
                                                                       loss 1.1484 (0.4657)
    epoch [12/100]
                   batch
                                                    data 0.000
                                                                                             lr 2.4261e-03 eta 0:03:41
    epoch [13/100] batch [20/20] time 0.096 (0.128) data 0.000 (0.017)
                                                                       loss 0.8467 (0.5009) lr 2.4122e-03 eta 0:03:42
    epoch [14/100]
                   batch [20/20]
                                 time 0.096 (0.128) data 0.000 (0.023) loss 0.5547
                                                                                    (0.4495) lr 2.3972e-03 eta 0:03:39
    epoch [15/100] batch [20/20] time 0.143 (0.152) data 0.000 (0.015) loss 1.0430 (0.5549) lr 2.3810e-03 eta 0:04:18
                                                    data 0.000 (0.032) loss 1.3906 (0.4799) lr 2.3638e-03 eta 0:05:48
    epoch [16/100]
                   batch [20/20] time 0.136 (0.207)
    epoch [17/100]
                   batch
                         [20/20]
                                 time 0.095 (0.129) data 0.000 (0.017)
                                                                       loss 0.0238
                                                                                    (0.3497)
                                                                                             lr 2.3454e-03 eta 0:03:34
    epoch [18/100]
                   batch
                         [20/20]
                                 time 0.110 (0.130)
                                                    data 0.000 (0.020) loss 0.1337
                                                                                    (0.2804) lr 2.3259e-03 eta 0:03:33
    epoch [19/100]
                  batch
                         [20/20]
                                 time 0.094 (0.130) data 0.000
                                                               (0.015)
                                                                       loss 1.0420
                                                                                    (0.3864) lr 2.3054e-03 eta 0:03:30
    epoch [20/100] batch [20/20] time 0.156 (0.146) data 0.000 (0.024) loss 0.3484 (0.4984) lr 2.2839e-03 eta 0:03:52
    epoch [21/100] batch [20/20]
                                time 0.141 (0.193) data 0.000
                                                               (0.036) loss 0.8184 (0.3434) lr 2.2613e-03 eta 0:05:05
    epoch [22/100] batch [20/20] time 0.096 (0.129) data 0.000 (0.018) loss 0.2090 (0.4361) lr 2.2377e-03 eta 0:03:21
1 # Accuracy on the New Classes.
2 args.model_dir = "outputs/cocoop"
3 args.output_dir = "outputs/cocoop/new_classes"
4 args.subsample_classes = "new"
5 args.load_epoch = 100
6 args.eval_only = True
7 coop_novel_acc = main(args)
→ Loading trainer: CoCoOp
    Loading dataset: EuroSAT
    Reading split from /content/ProMetaR/data/eurosat/split_zhou_EuroSAT.json
    Loading preprocessed few-shot data from /content/ProMetaR/data/eurosat/split fewshot/shot 16-seed 1.pkl
    SUBSAMPLE NEW CLASSES!
    Building transform_train
    + random resized crop (size=(224, 224), scale=(0.08, 1.0))
    + random flip
    + to torch tensor of range [0, 1]
    + normalization (mean=[0.48145466, 0.4578275, 0.40821073], std=[0.26862954, 0.26130258, 0.27577711])
    Building transform_test
    + resize the smaller edge to 224
    + 224x224 center crop
    + to torch tensor of range [0, 1]
    + normalization (mean=[0.48145466, 0.4578275, 0.40821073], std=[0.26862954, 0.26130258, 0.27577711])
   Dataset
               EuroSAT
    # classes
               5
    # train_x
               80
    # val
               20
               3,900
    # test
    Loading CLIP (backbone: ViT-B/16)
    /usr/local/lib/python3.10/dist-packages/torch/utils/data/dataloader.py:617: UserWarning: This DataLoader will create 8 w
      warnings.warn(
    /usr/local/lib/python3.10/dist-packages/torch/optim/lr_scheduler.py:62: UserWarning: The verbose parameter is deprecated
      warnings.warn(
    /content/ProMetaR/dassl/utils/torchtools.py:102: FutureWarning: You are using `torch.load` with `weights_only=False` (th
      checkpoint = torch.load(fpath, map_location=map_location)
    Building custom CLIP
    Initial context: "a photo of a"
    Number of context words (tokens): 4
    Turning off gradients in both the image and the text encoder
    Parameters to be updated: {'prompt_learner.ctx', 'prompt_learner.meta_net.linear2.weight', 'prompt_learner.meta_net.line
    Loading evaluator: Classification
    Loading weights to prompt_learner from "outputs/cocoop/prompt_learner/model.pth.tar-100" (epoch = 100)
    Evaluate on the *test* set
                  | 39/39 [01:01<00:00, 1.57s/it]=> result
    100%
    * total: 3,900
    * correct: 1,687
    * accuracy: 43.3%
    * error: 56.7%
    * macro_f1: 39.0%
```

Q3. Analyzing the results of CoCoOp

Compare the results of CoCoOp with those of CoOp that we trained in Lab Session 4. Discuss possible reasons for the performance differences observed between CoCoOp and CoOp.

The results show that CoCoOp performs well on the base classes with an accuracy of 90.8%, which is strong. However, on new classes, its performance drops dramatically to 43.3%. Compared to CoOp, CoOp has lower accuracy on new classes. This demonstrates CoOp's fixed prompt tuning generalizes better to unseen categories in this case.

Possible reasons for the performance difference could be that CoCoOp relies on dynamic, image-conditioned prompts and this causes overfitting to visual features of base classes. Thus, its ability to handle unseen categories. In contrast CoOp has a better generalization across



class distributions.