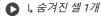
Homework 4

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.

```
import torch.nn as nn
class CoCoOpPromptLearner(nn.Module):
    def __init__(self, cfg, classnames, clip_model):
       super().__init__()
       n_cls = len(classnames)
       n_ctx = cfg.TRAINER.COCOOP.N_CTX
       ctx_init = cfg.TRAINER.COCOOP.CTX_INIT
       dtype = clip model.dtype
       ctx_dim = clip_model.In_final.weight.shape[0]
       vis_dim = clip_model.visual.output_dim
       clip_imsize = clip_model.visual.input_resolution
       cfg_imsize = cfg.INPUT.SIZE[0]
       assert cfg_imsize == clip_imsize, f"cfg_imsize ({cfg_imsize}) must equal to clip_imsize ({clip_imsize})"
        if ctx_init:
           # use given words to initialize context vectors
           ctx_init = ctx_init.replace("_", " ")
           n_ctx = len(ctx_init.split(" "))
           prompt = clip.tokenize(ctx_init)
           with torch.no_grad():
               embedding = clip_model.token_embedding(prompt).type(dtype)
           ctx_vectors = embedding[0, 1: 1 + n_ctx, :]
           prompt_prefix = ctx_init
           # random initialization
           ctx_vectors = torch.empty(n_ctx, ctx_dim, dtype=dtype)
           nn.init.normal_(ctx_vectors, std=0.02)
           prompt_prefix = " ".join(["X"] * n_ctx)
       print(f'Initial context: "{prompt_prefix}"')
       print(f"Number of context words (tokens): {n_ctx}")
        self.ctx = nn.Parameter(ctx_vectors) # Wrap the initialized prompts above as parameters to make them trainable.
        ### Tokenize ###
       classnames = [name.replace("_", " ") for name in classnames] # @|) "Forest"
       name_lens = [len(_tokenizer.encode(name)) for name in classnames]
       prompts = [prompt_prefix + " " + name + "." for name in classnames] # 예) "A photo of Forest."
```

tokenized_prompts = torch.cat([clip.tokenize(p) for p in prompts]) # 예) [49406, 320, 1125, 539...]

```
###### Q1. Fill in the blank ######
   ######## Define Meta Net ########
   self.meta_net = nn.Sequential(OrderedDict([
       ("linear1", nn.Linear(vis_dim, vis_dim // 16)),
       ("relu", nn.ReLU(inplace=True)),
       ("linear2", nn.Linear(vis_dim // 16, ctx_dim))
   1))
   ## Hint: meta network is composed to linear layer, relu activation, and linear layer.
   if cfg.TRAINER.COCOOP.PREC == "fp16":
      self.meta_net.half()
   with torch.no_grad():
       embedding = clip_model.token_embedding(tokenized_prompts).type(dtype)
   # These token vectors will be saved when in save_model(),
   # but they should be ignored in load_model() as we want to use
   # those computed using the current class names
   self.register_buffer("token_prefix", embedding[:, :1, :]) # SOS
   {\tt self.register\_buffer("token\_suffix",\ embedding[:,\ 1+n\_ctx:,\ :])} \ \ \#\ {\tt CLS},\ {\tt EOS}
   self.n_cls = n_cls
   self.n\_ctx = n\_ctx
   self.tokenized_prompts = tokenized_prompts # torch.Tensor
   self.name_lens = name_lens
def construct_prompts(self, ctx, prefix, suffix, label=None):
   # dimO is either batch_size (during training) or n_cls (during testing)
   # ctx: context tokens, with shape of (dim0, n_ctx, ctx_dim)
   # prefix: the sos token, with shape of (n_cls, 1, ctx_dim)
   # suffix: remaining tokens, with shape of (n_cls, *, ctx_dim)
   if label is not None:
      prefix = prefix[label]
      suffix = suffix[label]
   prompts = torch.cat(
      [
          prefix, # (dimO, 1, dim)
          ctx, # (dimO, n_ctx, dim)
          suffix, # (dimO, *, dim)
      dim=1,
   )
   return prompts
def forward(self, im_features):
   prefix = self.token_prefix
   suffix = self.token_suffix
   ctx = self.ctx # (n_ctx, ctx_dim)
   ######## Q2,3. Fill in the blank #######
   bias = self.meta_net(im_features) # (batch, ctx_dim)
   bias = bias.unsqueeze(1) # (batch, 1, ctx_dim)
   ctx = ctx.unsqueeze(0) # (1, n_ctx, ctx_dim)
   ctx_shifted = ctx + bias # (batch, n_ctx, ctx_dim)
```

Use instance-conditioned context tokens for all classes

```
prompts = []
       for ctx_shifted_i in ctx_shifted:
          ctx_i = ctx_shifted_i.unsqueeze(0).expand(self.n_cls, -1, -1)
          pts_i = self.construct_prompts(ctx_i, prefix, suffix) # (n_cls, n_tkn, ctx_dim)
          prompts.append(pts_i)
       prompts = torch.stack(prompts)
       return prompts
class CoCoOpCustomCLIP(nn.Module):
   def __init__(self, cfg, classnames, clip_model):
       super().__init__()
       self.prompt_learner = CoCoOpPromptLearner(cfg, classnames, clip_model)
       self.tokenized_prompts = self.prompt_learner.tokenized_prompts
       self.image_encoder = clip_model.visual
       self.text_encoder = TextEncoder(clip_model)
       self.logit_scale = clip_model.logit_scale
       self.dtype = clip_model.dtype
   def forward(self, image, label=None):
       tokenized prompts = self.tokenized prompts
       logit_scale = self.logit_scale.exp()
       image_features = self.image_encoder(image.type(self.dtype))
       image_features = image_features / image_features.norm(dim=-1, keepdim=True)
       ######## Q4. Fill in the blank #######
       prompts = self.prompt_learner(image_features)
       logits = []
       for pts_i, imf_i in zip(prompts, image_features):
           text_features = self.text_encoder(pts_i, tokenized_prompts)
           text_features = text_features / text_features.norm(dim=-1, keepdim=True)
           l_i = logit_scale * imf_i @ text_features.t()
           logits.append(l_i)
       logits = torch.stack(logits)
       if self.prompt_learner.training:
           return F.cross_entropy(logits, label)
       return logits
```

∨ 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.

```
# Train on the Base Classes Train split and evaluate accuracy on the Base Classes Test split.

args.trainer = "CoCoOp"

args.train_batch_size = 4

args.epoch = 100

args.output_dir = "outputs/cocoop"

args.subsample_classes = "base"

args.eval_only = False

cocoop_base_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 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 5
     # train_x
                80
     # val
                20
     # 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 worker pr
       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.linear1.weight', 'prompt_learner.meta_net.linear2.weig
     /usr/local/lib/python3.10/dist-packages/torch/optim/lr_scheduler.py:62: UserWarning: The verbose parameter is deprecated. Please
     /content/ProMetaR/dassI/utils/torchtools.py:102: FutureWarning: You are using `torch.load` with `weights_only=False` (the curren
       checkpoint = torch.load(fpath, map_location=map_location)
     Loading evaluator: Classification
     Found checkpoint at outputs/cocoop (will resume training)
     Loading checkpoint from "outputs/cocoop/prompt_learner/model.pth.tar-100"
     Loaded model weights
     Loaded optimizer
     Loaded scheduler
     Previous epoch: 100
     Finish training
     Deploy the last-epoch model
     Evaluate on the *test* set
     100%| 42/42 [01:07<00:00, 1.61s/it]=> result
     * total: 4.200
     * correct: 3,813
     * accuracy: 90.8%
     * error: 9.2%
     * macro_f1: 90.9%
     Elapsed: 0:01:07
# Accuracy on the New Classes.
args.model_dir = "outputs/cocoop"
args.output_dir = "outputs/cocoop/new_classes"
args.subsample_classes = "new'
args.load_epoch = 100
args.eval_only = True
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
     # test
                3,900
     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 worker pr
```

https://colab.research.google.com/drive/14xvIQWtbQCZ4GOduhn2gzcUMepb Bl8b?hl=ko#scrollTo=QeRABv42Ku4E&printMode=true

```
warnings.warn(
```

/usr/local/lib/python3.10/dist-packages/torch/optim/lr_scheduler.py:62: UserWarning: The verbose parameter is deprecated. Please warnings.warn(

/content/ProMetaR/dassl/utils/torchtools.py:102: FutureWarning: You are using `torch.load` with `weights_only=False` (the curren 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.linear1.weight', 'prompt_learner.meta_net.linear2.weig Loading evaluator: Classification

Loading weights to prompt_learner from "outputs/cocoop/prompt_learner/model.pth.tar-100" (epoch = 100)

Evaluate on the *test* set

100%| 39/39 [00:59<00:00, 1.51s/it]=> result

* 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.

CoOp 결과는 다음과 같다.

base class:

accuracy: 91.4%

• error: 8.6%

• macro_f1: 91.5%

new class:

accuracy: 51.5%

error: 48.5%

macro_f1: 45.6%

CoCoOp 결과는 다음과 같다.

base class:

• accuracy: 90.8%

• error: 9.2%

• macro_f1: 90.9%

new class:

accuracy: 43.3%

• error: 56.7%

macro_f1: 39.0%

CoOp의 경우 base class에 대해 학습된 prompt를 그대로 new class에 적용한다.

그러나 CoCoOp의 경우 input에 대해 meta network를 적용해 그 output을 prompt에 더함으로써, 새로운 input에 대해 더 알맞은 prompt를 산출해 new class에 적용한다.

따라서 이론적으로는 CoCoOp이 CoOp에 대해 base class 성능이 비슷하고 new class에서의 성능은 더 좋아야 한다.

그러나 실제로는 CoCoOp이 CoOp에 비해 base class 뿐 아니라 특히 new class에서 성능이 확연히 뒤떨어졌다. 그 이유는 다음과 같은 것으로 예상된다.

class 수가 적음 (overfitting):

적은 class 수 (5개) 로 인해 new class 실험에서의 training set에 대해 meta network가 overfitting 되어, test accuracy가 낮아졌을 가 능성이 있다.

base class의 경우 base class에 맞는 prompt를 활용했으므로 meta network의 overfitting 에 의한 오류는 상대적으로 영향이 적었을 것이다.

코딩을 시작하거나 AI로 코드를 <u>생성</u>하세요.