

```
1 import os
2 from datetime import datetime
3
4 # For this notebook to run with updated APIs, we need torch 1.12+ and torchvision
5 # 0.13+
6 try:
7     import torch
8     import torchvision
9     assert int(torch.__version__.split(".")[1]) >= 12 or
10    int(torch.__version__.split(".")[0]) == 2, "torch version should be 1.12+"
11    assert int(torchvision.__version__.split(".")[1]) >= 13, "torchvision version
12    should be 0.13+"
13    print(f"torch version: {torch.__version__}")
14    print(f"torchvision version: {torchvision.__version__}")
15 except:
16     print("[INFO] torch/torchvision versions not as required, installing nightly
17     versions.")
18     #!pip3 install -U torch torchvision torchaudio --index-url
19     https://download.pytorch.org/whl/cu118
20
21     import torch
22     import torchvision
23     print(f"torch version: {torch.__version__}")
24     print(f"torchvision version: {torchvision.__version__}")
25
26
27 import os
28 import shutil
29 import random
30 import matplotlib.pyplot as plt
31 import torch
32 import torchvision
33 from torch import nn
34 from torchvision import transforms
35 try:
36     from torchinfo import summary
37 except:
38     print("[INFO] Couldn't find torchinfo... installing it.")
39     #!pip install -q torchinfo
40     from torchinfo import summary
41 from pathlib import Path
42 import os
43
44 test_run = 7
45 EPOCHS = 200
46
47 print("#####")
48 print("#####")
49 print("#####")
50 print(f"Vision Transformer Experiment For Test Run: {test_run}")
51 print(f"The Number of Epochs are: {EPOCHS}")
52 print("#####")
53 print("#####")
54 device = "cuda" if torch.cuda.is_available() else "cpu"
55 print("#####")
```

```

56 print(f"The device for running the experiment is {device}")
57 print("#####")
58
59
60
61
62
63
64 # Path to the dataset folder
65 main_folder = f"../../training_data_split_run_{test_run}"
66 splits = ['train', 'test', 'val']
67
68 summary = {}
69
70 for split in splits:
71     split_path = os.path.join(main_folder, split)
72     total_files = 0
73     summary[split] = {}
74
75
76     if not os.path.exists(split_path):
77         print(f"Directory not found: {split_path}")
78         continue
79
80     for disease in os.listdir(split_path):
81         disease_path = os.path.join(split_path, disease)
82
83         if not os.path.isdir(disease_path):
84             continue
85
86
87         num_files = len([
88             f for f in os.listdir(disease_path)
89             if os.path.isfile(os.path.join(disease_path, f))
90         ])
91
92         summary[split][disease] = num_files
93         total_files += num_files
94
95     print(f"\n{split.upper()} - Total files: {total_files}")
96     for disease, count in summary[split].items():
97         print(f"  {disease}: {count} files")
98
99
100 train_dir = f"../../training_data_split_run_{test_run}/train"
101 test_dir = f"../../training_data_split_run_{test_run}/val"
102
103 # Create image size (from Table 3 in the ViT paper)
104 IMG_SIZE = 224
105
106 # Create transform pipeline manually
107 manual_transforms = transforms.Compose([
108     transforms.Resize((IMG_SIZE, IMG_SIZE)),
109     transforms.ToTensor(),
110 ])
111 print(f"Manually created transforms: {manual_transforms}")
112
113 from going_modular import data_setup
114
```

```
115 BATCH_SIZE = 32 # this is lower than the ViT paper but it's because we're starting
116 # small
117
118 # Create data loaders
119 train_dataloader, test_dataloader, class_names = data_setup.create_dataloaders(
120     train_dir=train_dir,
121     test_dir=test_dir,
122     transform=manual_transforms, # use manually created transforms
123     batch_size=BATCH_SIZE
124 )
125
126 train_dataloader, test_dataloader, class_names
127
128 print(f"The class names are: {class_names}")
129
130
131 from going_modular.helper_functions import set_seeds
132 set_seeds()
133
134
135 pretrained_vit_weights = torchvision.models.ViT_B_16_Weights.DEFAULT
136
137 # 2. Setup a ViT model instance with pretrained weights
138 pretrained_vit =
139     torchvision.models.vit_b_16(weights=pretrained_vit_weights).to(device)
140
141 # 3. Freeze the base parameters
142 #for parameter in pretrained_vit.parameters():
143 #    parameter.requires_grad = False
144
145 # 4. Change the classifier head (set the seeds to ensure same initialization with
146 # linear head)
147 #set_seeds()
148 pretrained_vit.heads = nn.Linear(in_features=768,
149                                 out_features=len(class_names)).to(device)
150 # pretrained_vit
151
152
153 from torchinfo import summary
154
155 # Print a summary using torchinfo (uncomment for actual output)
156 summary(model=pretrained_vit,
157           input_size=(32, 3, 224, 224), # (batch_size, color_channels, height, width)
158           # col_names=["input_size"], # uncomment for smaller output
159           col_names=["input_size", "output_size", "num_params", "trainable"],
160           col_width=20,
161           row_settings=["var_names"])
162
163 # Get automatic transforms from pretrained ViT weights
164 pretrained_vit_transforms = pretrained_vit_weights.transforms()
165 print(pretrained_vit_transforms)
166
167 train_dataloader_pretrained, test_dataloader_pretrained, class_names =
168     data_setup.create_dataloaders(train_dir=train_dir,
169                                   test_dir=test_dir,
```

```
168         transform=pretrained_vit_transforms,
169         batch_size=32)
170
171 from going_modular import engine
172
173 # Create optimizer and loss function
174 #optimizer = torch.optim.Adam(params=pretrained_vit.parameters(),
175 #                             lr=1e-3)
176 optimizer = torch.optim.Adam(params=pretrained_vit.parameters(),
177                             lr=1e-5)
178
179
180 loss_fn = torch.nn.CrossEntropyLoss()
181
182
183 import torch
184
185 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
186
187 # Wrap model for multi-GPU if available
188 if torch.cuda.device_count() > 1:
189     print(f"[INFO] Using DataParallel on {torch.cuda.device_count()} GPUs")
190     #pretrained_vit = torch.nn.DataParallel(pretrained_vit)
191     pretrained_vit = torch.nn.DataParallel(pretrained_vit)
192
193 pretrained_vit.to(device)
194
195
196 # Train the classifier head of the pretrained ViT feature extractor model
197 #set_seeds()
198
199 from going_modular.helper_functions import plot_loss_curves
200 import os
201 from datetime import datetime
202 import numpy
203 import pickle
204 import numpy as np
205
206
207 timestamp = datetime.now().strftime("%Y-%m-%d_%H-%M-%S")
208 results_directory =
209 f"Results_Vision_Transformer_experiment_run_{test_run}_{timestamp}"
210 os.makedirs(results_directory, exist_ok=True)
211
212 import time
213 start_time = time.time()
214 pretrained_vit_results = engine.train(model=pretrained_vit,
215                                         train_dataloader=train_dataloader_pretrained,
216                                         test_dataloader=test_dataloader_pretrained,
217                                         optimizer=optimizer,
218                                         loss_fn=loss_fn,
219                                         epochs=EPOCHS,
220                                         device=device,
221                                         test_run = test_run,
222                                         results_directory=results_directory
223                                         )
224 training_time = time.time() - start_time
225 print(f"The total time for training is {training_time:.2f} seconds")
```

```

225
226
227
228
229 plot_loss_curves(pretrained_vit_results,
230     save_path=results_directory+f"/loss_accuracy_curve_{test_run}_{timestamp}.png")
231
232 # Save the model
233 from going_modular import utils
234
235 utils.save_model(model=pretrained_vit,
236                     target_dir=results_directory,
237                     model_name=f"Pretrained_vit_scratch_training_test_run_{test_run}_save_after_200_epochs.pth")
238
239
240 from torchinfo import summary
241 # Print a summary using torchinfo (uncomment for actual output)
242 summary(model=pretrained_vit,
243           input_size=(32, 3, 224, 224), # (batch_size, color_channels, height, width)
244           # col_names=["input_size"], # uncomment for smaller output
245           col_names=["input_size", "output_size", "num_params", "trainable"],
246           col_width=20,
247           row_settings=["var_names"])
248 )
249
250
251 from pathlib import Path
252
253 # Get the model size in bytes then convert to megabytes
254 pretrained_vit_model_size =
255     Path(results_directory+f"/Pretrained_vit_scratch_training_test_run_{test_run}_save_after_200_epochs.pth").stat().st_size // (1024*1024) # division converts bytes to
256     megabytes (roughly)
257 print(f"Pretrained ViT feature extractor model size: {pretrained_vit_model_size} MB")
258
259 # Setup directory paths to train and test images
260 val_dir =
261     f"/leonardo_work/IscrC_ArtLLMs/alitariqnagi_work/training_data_split_run_{test_run}/val"
262 test_dir =
263     f"/leonardo_work/IscrC_ArtLLMs/alitariqnagi_work/training_data_split_run_{test_run}/test"
264
265 from going_modular import data_setup
266 from going_modular.data_setup import create_dataloaders_no_shuffle
267 # Setup dataloaders create_dataloaders_no_shuffle
268 val_dataloader_pretrained, test_dataloader_pretrained, class_names =
269     data_setup.create_dataloaders_no_shuffle(train_dir=val_dir,
270
271             test_dir=test_dir,
272
273             transform=pretrained_vit_transforms,
274
275             batch_size=32)

```

```

271
272 ##########
273 import numpy as np
274 import torch
275 import torch.nn.functional as F
276 from torch.nn.functional import softmax
277 import numpy
278 from tqdm import tqdm
279
280 alpha = 0.1
281
282
283 #####-----Calculation of
284 Calibration Scores-----#
285
286 def calibration_scores(model, data_loader, device):
287     model.eval()
288     scores_calibration_set = []
289
290     with torch.no_grad():
291         for X, y in tqdm(data_loader, desc="#Proceeding with the
292 Calculation of Calibration Scores#####"):
293             X = X.to(device)
294             y = y.to(device)
295
296             y = y.cpu().numpy()
297
298             pred_logits = model(X)
299
300             # cal_smx = model(calib_X).softmax(dim=1 ).numpy()
301             cal_smx = softmax(pred_logits, dim=1).cpu().numpy()
302
303             # cal_scores = 1-cal_smx[np . arange(n),cal_labels]
304             scores_per_batch = 1 - cal_smx[np.arange(y.shape[0])], y
305
306             scores_calibration_set.extend(scores_per_batch)
307
308     scores_calibration_set = numpy.array(scores_calibration_set)
309
310     return scores_calibration_set
311 #####
312
313 #####-----Calculation of
314 Prediction Sets-----#
315
316 def predict_with_conformal_sets(model, data_loader, calibration_scores_array, alpha,
317 device):
318
319     model.eval()
320
321     prediction_sets = []
322     true_labels = []
323
324     n = len(calibration_scores_array)

```

```

324
325     #get adjusted quantile
326     #q_level = np . ceil((n+1 ) *(1-alpha)) /n
327     q_level = np.ceil((n + 1) * (1 - alpha)) / n
328
329     # qhat = np . quantile(calibration_scores_array, q_level, method='higher')
330     qhat = np.quantile(calibration_scores_array, q_level, method='higher')
331
332     with torch.no_grad():
333         for X, y in tqdm(data_loader, desc="# ##### Proceeding with the
Calculation of Testing Scores#####"):
334
335             X = X.to(device)
336             y = y.cpu().numpy()
337             pred_logits = model(X)
338
339             # cal_smx = model(calib_X).softmax(dim=1 ).numpy()
340             val_smx = softmax(pred_logits, dim=1)
341             val_smx = val_smx.cpu().numpy()
342
343             # prediction_sets = val_smx >= (1-qhat)
344             masks = val_smx >= (1 - qhat)
345
346             for _ in masks:
347                 prediction_sets.append(np.where(_)[0].tolist())
348
349             true_labels.extend(y.tolist())
350
351     return prediction_sets, true_labels
352 #####-----EVALUATION-----
353
354
355 #####----- Run Full Conformal Prediction -----
356 def evaluate(prediction_sets, labels):
357     number_of_correct_predictions = 0
358     for pred_set, true_label in zip(prediction_sets, labels):
359         if true_label in pred_set:
360             number_of_correct_predictions += 1
361     empirical_test_coverage = number_of_correct_predictions / len(labels)
362     average_prediction_set_size = np.mean([len(s) for s in prediction_sets])
363     return empirical_test_coverage, average_prediction_set_size
364 #####-----#
365
366 scores_calibration = calibration_scores(model=pretrained_vit,
367 data_loader=val_dataloader_pretrained, device=device)
368 prediction_sets, true_labels = predict_with_conformal_sets(model=pretrained_vit,
369 data_loader=test_dataloader_pretrained, calibration_scores_array=scores_calibration,
\                                         alpha=alpha,
370                                         device=device)
371 test_set_coverage, avg_size = evaluate(prediction_sets=prediction_sets,
372 labels=true_labels)
373 import os
374 from datetime import datetime
375 import numpy

```

```

375 import pickle
376 import numpy as np
377
378
379 timestamp = datetime.now().strftime("%Y-%m-%d_%H-%M-%S")
380 os.makedirs(f"Results_normal_Conformal_Prediction_experiment_run_{test_run}_{timestamp}",
381 p)", exist_ok=True)
382 results_directory =
383 f"Results_normal_Conformal_Prediction_experiment_run_{test_run}_{timestamp}"
384
385 import pickle
386
387 with open(f"{results_directory}/scores_calibration_run_{test_run}_{timestamp}.pkl",
388 "wb") as f:
389     pickle.dump(scores_calibration, f)
390
391 # with open("scores_calibration.pkl", "rb") as f:
392 #     scores_calibration = pickle.load(f)
393
394
395 with
396     open(f"Results_normal_Conformal_Prediction_experiment_run_{test_run}_{timestamp}/logs_
397 _experiment_run_{test_run}_{timestamp}.txt", "w") as f:
398     f.write("##### Results
399 #####\n")
400     f.write(f"Coverage: {test_set_coverage:.3f}\n")
401     f.write(f"Average Prediction Set Size for Test Set: {avg_size:.3f}\n")
402
403     f.write("#####\n")
404 print(f"logs file in directory
405 Results_normal_Conformal_Prediction_experiment_run_{test_run}_{timestamp}")
406
407
408 print("#####Results:#####")
409 print(f"Coverage: {test_set_coverage:.3f}")
410 print(f"Average Prediction Set Size for Test Set: {avg_size:.3f}")
411 print("#####")
412
413 np.save(f"Results_normal_Conformal_Prediction_experiment_run_{test_run}_{timestamp}/c_
414 alibration_scores_experiment_run_{test_run}_{timestamp}.npy", scores_calibration)
415
416
417
418 #####TO BE
419 CHANGED#####
420 #####
421 import numpy as np
422 import matplotlib.pyplot as plt
423 from collections import Counter

```

```

422 import pandas as pd
423
424 set_sizes = [len(prediction_set) for prediction_set in prediction_sets]
425 counts = Counter(set_sizes)
426 total = sum(counts.values())
427 sizes = sorted(counts.keys())
428 frequencies = [counts[size] for size in sizes]
429 percentages = [100 * freq / total for freq in frequencies]
430
431 data = [{'size': len(prediction_set), 'correct': int(true in prediction_set)} for
432 prediction_set, true in zip(prediction_sets, true_labels)]
433 df = pd.DataFrame(data)
434 coverage_stats = df.groupby('size').agg({'correct': ['mean', 'sum', 'count']})
435 coverage_stats.columns = ['coverage', 'correct_count', 'total_count']
436 coverage_stats = coverage_stats.reindex(sizes, fill_value=0)
437
438 cumulative_correct = coverage_stats['correct_count'].cumsum()
439 cumulative_coverage = 100 * cumulative_correct / len(true_labels)
440 max_freq = max(frequencies)
441 cumulative_scaled = (cumulative_coverage / 100) * max_freq
442
443 fig, ax = plt.subplots(figsize=(13, 6))
444 bars = ax.bar(sizes, frequencies, color='pink', edgecolor='black', label='Frequency',
445 alpha=0.7)
446 ax.plot(sizes, cumulative_scaled, color='red', marker='o', linestyle='--', alpha=0.4,
447 linewidth=2, label='Cumulative Coverage')
448
449 for bar, size, pct in zip(bars, sizes, percentages):
450     height = bar.get_height()
451     count = counts[size]
452     coverage = coverage_stats.loc[size, 'coverage'] * 100
453     ax.annotate(f'{count} ({pct:.1f}%) \n Cov: {coverage:.1f}%', xy=(bar.get_x() + bar.get_width() / 2, height / 2),
454         ha='center', va='center', fontsize=9)
455
456 for bar, cum_cov in zip(bars, cumulative_coverage):
457     height = bar.get_height()
458     ax.annotate(f'CumCov: {cum_cov:.1f}%',
459                 xy=(bar.get_x() + bar.get_width() / 2, height),
460                 xytext=(0, 8),
461                 textcoords="offset points",
462                 ha='center', va='bottom', fontsize=9, color='darkorange')
463
464 mean_size = np.mean(set_sizes)
465 median_size = np.median(set_sizes)
466 ax.axvline(mean_size, color='gray', linestyle='--', label=f'Mean = {mean_size:.2f}')
467 ax.axvline(median_size, color='orange', linestyle='--', label=f'Median = {median_size:.2f}')
468
469 ax.set_xlabel("Prediction Set Size")
470 ax.set_ylabel("Frequency")
471 ax.set_title("Prediction Set Size Distribution with Coverage (Test Set)")
472 ax.set_xticks(sizes)
473 ax.grid(axis='y', linestyle='--', alpha=0.6)
474 ax.legend(loc='upper right')
475 plt.tight_layout()
476 plt.savefig(os.path.join(results_directory,
477                         f"Prediction_Set_Size_Distribution_with_Coverage_(Test_Set)_{{test_run}}_{{timestamp}}.p
478 ng"))

```

```

476
477
478
479
480 plt.show()
481
482
483
484 #####script#####
485 from collections import Counter
486 import matplotlib.pyplot as plt
487 from collections import defaultdict
488 import numpy as np
489 from scipy.stats import pearsonr, spearmanr
490
491
492 def feature_stratified_coverage_by_class(prediction_sets, true_labels, alpha=0.1,
493 class_names=None):
494
495     class_to_indices = defaultdict(list)
496     for idx, label in enumerate(true_labels):
497         class_to_indices[label].append(idx)
498
499     class_coverages = {}
500     for cls, indices in class_to_indices.items():
501         covered = sum(true_labels[i] in prediction_sets[i] for i in indices)
502         coverage = covered / len(indices)
503         class_coverages[cls] = coverage
504
505     print("\n Feature-Stratified Coverage (by True Class Label):")
506     for cls, cov in sorted(class_coverages.items()):
507         name = class_names[cls] if class_names else str(cls)
508         print(f" Class {name:<20}: Coverage = {cov:.3f} (Target ≥ {1 - alpha:.2f})")
509
510     fsc_metric = min(class_coverages.values())
511     print(f"\n FSC Metric (minimum class-wise coverage): {fsc_metric:.3f}")
512
513     return fsc_metric, class_coverages
514
515
516 fsc_metric, per_class_cov = feature_stratified_coverage_by_class(
517     prediction_sets,
518     true_labels,
519     alpha=alpha,
520     class_names=class_names
521 )
522
523
524 def plot_class_wise_coverage(class_coverages, alpha=0.1, class_names=None,
525 title="Class-wise Coverage (FSC)"):
526
527     classes = list(class_coverages.keys())
528     coverages = [class_coverages[c] for c in classes]
529
530     # Get display names
531     labels = [class_names[c] if class_names else str(c) for c in classes]
532
533     # Plot

```

```

533 plt.figure(figsize=(12, 6))
534 bars = plt.bar(labels, coverages, color="pink", edgecolor="black")
535
536 # Target line
537 target = 1 - alpha
538 plt.axhline(target, color='red', linestyle='--', label=f"Target Coverage ({target:.2f})")
539
540 # Highlight underperforming classes
541 for bar, cov in zip(bars, coverages):
542     if cov < target:
543         bar.set_color('salmon')
544
545 plt.xticks(rotation=45, ha="right")
546 plt.ylim(0, 1.05)
547 plt.ylabel("Coverage")
548 plt.title(title)
549 plt.legend()
550 plt.tight_layout()
551 filename = f"Class_wise_Coverage_(FSC)_{test_run}_{timestamp}.png"
552 plt.savefig(os.path.join(results_directory, filename))
553 plt.show()
554
555 fsc_metric, class_coverages = feature_stratified_coverage_by_class(
556     prediction_sets,
557     true_labels,
558     alpha=alpha,
559     class_names=class_names
560 )
561
562 plot_class_wise_coverage(class_coverages, alpha=alpha, class_names=class_names)
563
564
565
566
567
568 def plot_coverage_vs_class_frequency(class_coverages, true_labels, alpha=0.1,
569                                         class_names=None, title="Coverage vs Class Frequency"):
570
571
572     class_indices = sorted(class_coverages.keys())
573     coverages = [class_coverages[c] for c in class_indices]
574
575
576     total = len(true_labels)
577     counts = Counter(true_labels)
578     frequencies = [counts[c] / total for c in class_indices]
579
580
581     labels = [class_names[c] if class_names else str(c) for c in class_indices]
582     x = np.arange(len(class_indices))
583
584     fig, ax1 = plt.subplots(figsize=(12, 6))
585
586
587     bars1 = ax1.bar(x - 0.2, coverages, width=0.4, label="Coverage", color='pink',
588                     edgecolor='black')
589     ax1.axhline(1 - alpha, color='red', linestyle='--', label=f"Target Coverage ({1 - alpha:.2f})")

```

```

589     ax1.set_ylim(0, 1.05)
590     ax1.set_ylabel("Coverage", color='pink')
591     ax1.tick_params(axis='y', labelcolor='pink')
592
593 # Highlight under-covered bars
594 for bar, cov in zip(bars1, coverages):
595     if cov < (1 - alpha):
596         bar.set_color('salmon')
597
598 # Twin y-axis for frequencies
599 ax2 = ax1.twinx()
600 bars2 = ax2.bar(x + 0.2, frequencies, width=0.4, label="Class Frequency",
color='gray', alpha=0.5)
601 ax2.set_ylabel("Class Frequency", color='gray')
602 ax2.tick_params(axis='y', labelcolor='gray')
603
604 # X-axis
605 plt.xticks(x, labels, rotation=45, ha="right")
606
607 # Title and legend
608 plt.title(title)
609 fig.legend(loc='upper right', bbox_to_anchor=(0.85, 0.85))
610 plt.tight_layout()
611 filename = f"plot_coverage_vs_class_frequency_{test_run}_{timestamp}.png"
612 plt.savefig(os.path.join(results_directory, filename))
613 plt.show()
614
615
616 # Step 1: Compute coverage by class
617 fsc_metric, class_coverages = feature_stratified_coverage_by_class(
618     prediction_sets,
619     true_labels,
620     alpha=alpha,
621     class_names=class_names
622 )
623
624 # Step 2: Compare against class frequency
625 plot_coverage_vs_class_frequency(class_coverages, true_labels, alpha=alpha,
class_names=class_names)
626
627
628 def plot_coverage_vs_class_frequency_with_correlation(class_coverages, true_labels,
alpha=0.1, class_names=None, title="Coverage vs Class Frequency"):
629
630
631
632     class_indices = sorted(class_coverages.keys())
633     coverages = [class_coverages[c] for c in class_indices]
634
635 # Compute frequencies
636     total = len(true_labels)
637     counts = Counter(true_labels)
638     frequencies = [counts[c] / total for c in class_indices]
639
640 # Class labels
641     labels = [class_names[c] if class_names else str(c) for c in class_indices]
642     x = np.arange(len(class_indices))
643
644 # Correlation analysis
645     pearson_corr, pearson_p = pearsonr(frequencies, coverages)

```

```

646     spearman_corr, spearman_p = spearmanr(frequencies, coverages)
647
648     print(f"\n Correlation between class frequency and coverage:")
649     print(f" - Pearson r = {pearson_corr:.3f}, p = {pearson_p:.3e}")
650     print(f" - Spearman p = {spearman_corr:.3f}, p = {spearman_p:.3e}")
651
652     # Plotting
653     fig, ax1 = plt.subplots(figsize=(12, 6))
654
655     bars1 = ax1.bar(x - 0.2, coverages, width=0.4, label="Coverage", color='pink',
656                      edgecolor='black')
656     ax1.axhline(1 - alpha, color='red', linestyle='--', label=f"Target Coverage ({1 -
alpha:.2f})")
657     ax1.set_ylimits(0, 1.05)
658     ax1.set_ylabel("Coverage", color='pink')
659     ax1.tick_params(axis='y', labelcolor='pink')
660
661     # Highlight under-covered classes
662     for bar, cov in zip(bars1, coverages):
663         if cov < (1 - alpha):
664             bar.set_color('salmon')
665
666     # Frequency bars on second axis
667     ax2 = ax1.twinx()
668     bars2 = ax2.bar(x + 0.2, frequencies, width=0.4, label="Class Frequency",
669                      color='gray', alpha=0.5)
670     ax2.set_ylabel("Class Frequency", color='gray')
671     ax2.tick_params(axis='y', labelcolor='gray')
672
673     # X-axis and titles
674     plt.xticks(x, labels, rotation=45, ha="right")
675     plt.title(title)
676     fig.legend(loc='upper right', bbox_to_anchor=(0.85, 0.85))
677     plt.tight_layout()
678     filename =
679     f"plot_coverage_vs_class_frequency_with_correlation_{test_run}_{timestamp}.png"
680     plt.savefig(os.path.join(results_directory, filename))
681     plt.show()
682
683 plot_coverage_vs_class_frequency_with_correlation(
684     class_coverages=class_coverages,
685     true_labels=true_labels,
686     alpha=alpha,
687     class_names=class_names
688 )
689
690
691 def size_stratified_coverage(prediction_sets, true_labels, alpha=0.1, bins=[[1], [2],
692 [3, 4, 5], list(range(6, 100))]):
693     """
694         Computes size-stratified coverage (SSC): coverage within prediction set size
695         bins.
696
697     Args:
698         prediction_sets (List[List[int]]): output sets from conformal_predict
699         true_labels (List[int]): ground-truth labels
700         alpha (float): target error level (e.g., 0.1)
701         bins (List[List[int]]): list of prediction set size groups

```

```

700
701     Returns:
702         ssc (float): min group coverage (worst-case)
703         bin_coverages (dict): {bin_range: coverage}
704     """
705     bin_groups = defaultdict(list)
706
707     # Assign examples to bins based on set size
708     for i, pred_set in enumerate(prediction_sets):
709         set_size = len(pred_set)
710         for bin_range in bins:
711             if set_size in bin_range:
712                 bin_groups[str(bin_range)].append(i)
713                 break
714
715     bin_coverages = {}
716     for bin_key, indices in bin_groups.items():
717         covered = sum(true_labels[i] in prediction_sets[i] for i in indices)
718         coverage = covered / len(indices) if indices else 0
719         bin_coverages[bin_key] = coverage
720
721     print("\n Size-Stratified Coverage (SSC):")
722     for bin_key, cov in bin_coverages.items():
723         print(f" Set Size {bin_key}:10s: Coverage = {cov:.3f} (Target ≥ {1 - alpha:.2f})")
724
725     ssc = min(bin_coverages.values())
726     print(f"\n SSC Metric (min bin-wise coverage): {ssc:.3f}")
727     return ssc, bin_coverages
728
729
730 ssc_metric, bin_coverages = size_stratified_coverage(
731     prediction_sets,
732     true_labels,
733     alpha=alpha,
734     bins=[[1], [2], [3], list(range(4, 100))])
735 )
736
737
738
739 def plot_ssc_coverage(bin_coverages, alpha=0.1, title="Size-Stratified Coverage"):
740     """
741         Plots coverage per prediction set size bin.
742
743         Args:
744             bin_coverages (dict): mapping from bin label (e.g. '[1]', '[2, 3]') to
745             coverage
746                 alpha (float): target error level
747                 title (str): plot title
748             """
749             bin_labels = list(bin_coverages.keys())
750             coverages = [bin_coverages[k] for k in bin_labels]
751
752             plt.figure(figsize=(10, 5))
753             bars = plt.bar(bin_labels, coverages, color='pink', edgecolor='black')
754
755             # Target line
756             target = 1 - alpha
757             plt.axhline(target, color='red', linestyle='--', label=f"Target Coverage
({target:.2f})")

```

```

757
758 # Highlight low-coverage bins
759 for bar, cov in zip(bars, coverages):
760     if cov < target:
761         bar.set_color('grey')
762
763 plt.xticks(rotation=45, ha='right')
764 plt.ylim(0, 1.05)
765 plt.ylabel("Coverage")
766 plt.title(title)
767 plt.legend()
768 plt.tight_layout()
769 filename = f"plot_ssc_coverage_{test_run}_{timestamp}.png"
770 plt.savefig(os.path.join(results_directory, filename))
771 plt.show()
772
773
774 ssc_metric, bin_coverages = size_stratified_coverage(
775     prediction_sets,
776     true_labels,
777     alpha=alpha,
778     bins=[[1], [2], [3], list(range(4, 100))])
779 )
780
781 plot_ssc_coverage(bin_coverages, alpha=alpha)
782
783
784
785 def plot_prediction_set_size_by_class(prediction_sets, true_labels, class_names=None,
786                                         results_dir="", test_run="", timestamp=""):
787
788     # Group prediction set sizes by true class
789     size_per_class = defaultdict(list)
790     for pred_set, true_label in zip(prediction_sets, true_labels):
791         size_per_class[true_label].append(len(pred_set))
792
793     # Sort classes
794     sorted_classes = sorted(size_per_class.keys())
795     means = [np.mean(size_per_class[c]) for c in sorted_classes]
796     stds = [np.std(size_per_class[c]) for c in sorted_classes]
797     labels = [class_names[c] if class_names else str(c) for c in sorted_classes]
798
799     # Plot
800     plt.figure(figsize=(12, 6))
801     plt.bar(labels, means, yerr=stds, capsize=5, color="lightgreen",
802             edgecolor="black")
803     plt.ylabel("Avg Prediction Set Size")
804     plt.xlabel("Class")
805     plt.title("Average Prediction Set Size by Class")
806     plt.xticks(rotation=45, ha="right")
807     plt.grid(axis="y", linestyle="--", alpha=0.5)
808     plt.tight_layout()
809
810     # Save
811     filename = f"prediction_set_size_by_class_{test_run}_{timestamp}.png"
812     plt.savefig(os.path.join(results_directory, filename))
813     plt.show()
814
815 plot_prediction_set_size_by_class(

```

```

815     prediction_sets=prediction_sets,
816     true_labels=true_labels,
817     class_names=class_names,
818     results_dir=results_directory,
819     test_run=test_run,
820     timestamp=timestamp
821 )
822
823
824
825
826
827 def plot_calibration_score_distribution(calibration_scores, results_dir="", test_run="", timestamp ""):
828
829     plt.figure(figsize=(10, 5))
830     plt.hist(calibration_scores, bins=30, color="pink", edgecolor="black", alpha=0.8)
831     plt.xlabel("Nonconformity Score (1 - P(True Class))")
832     plt.ylabel("Frequency")
833     plt.title("Calibration Score Distribution")
834     plt.grid(axis='y', linestyle='--', alpha=0.6)
835     plt.tight_layout()
836
837     filename = f"calibration_score_distribution_{test_run}_{timestamp}.png"
838     plt.savefig(os.path.join(results_directory, filename))
839
840     plt.show()
841
842
843 plot_calibration_score_distribution(
844     calibration_scores=scores_calibration,
845     results_dir=results_directory,
846     test_run=test_run,
847     timestamp=timestamp
848 )
849
850
851 import json
852
853 metrics = {
854     "test_run": test_run,
855     "timestamp": timestamp,
856     "alpha": alpha,
857     "test_set_coverage": test_set_coverage,
858     "average_prediction_set_size": avg_size,
859     "FSC": fsc_metric,
860     "SSC": ssc_metric,
861     "class_coverages": {str(k): float(v) for k, v in class_coverages.items()},
862     "bin_coverages": bin_coverages
863 }
864
865 with open(os.path.join(results_directory,
866     f"metrics_summary_{test_run}_{timestamp}.json"), "w") as f:
867     json.dump(metrics, f, indent=4)
868
869
870
871 #####
```

872

873

#####