

```
1 import numpy as np
2 import torch
3 from torch import nn
4 from torch import optim
5 import torch.nn.functional as F
6 from torchvision import datasets, transforms, models
7 import torchvision
8 from torch.utils.data.sampler import SubsetRandomSampler
9 import matplotlib.pyplot as plt
10 import imageio
11 from torch.optim import lr_scheduler
12 from torchsummary import summary
13
14 # Device configuration
15 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
16 device
17
18 test_run = 7
19
20 print("test_run",test_run)
21
22
23 # Specify transforms using torchvision.transforms as transforms library
24 transformations = transforms.Compose([
25     transforms.Resize(255),
26     transforms.CenterCrop(224),
27     transforms.ToTensor(),
28     transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
29 ])
30
31 # Split the data into training, validation and datasets
32 # import splitfolders
33 # import os
34
35 # path = 'F:\skin_disease_data' #'mix-data'
36 # print(os.listdir(path))
37 # splitfolders.ratio(path, output="mix_data_split", seed=1337, ratio=(.6, .2, .2))
38
39 train_set_mix =
40     datasets.ImageFolder(f"../../training_data_split_run_{test_run}/train",
41     transform=transformations)
42 val_set_mix = datasets.ImageFolder(f"../../training_data_split_run_{test_run}/val",
43     transform=transformations)
44 test_set_mix =
45     datasets.ImageFolder(f"../../training_data_split_run_{test_run}/test",
46     transform=transformations)
47
48
49
50 class_names = train_set_mix.classes
51 class_names
52
53 train_loader_mix = torch.utils.data.DataLoader(train_set_mix, batch_size=32,
54     shuffle=True, num_workers=8)
```

```

54 val_loader_mix = torch.utils.data.DataLoader(val_set_mix, batch_size =32,
55 shuffle=False, num_workers=8)
56 test_loader_mix = torch.utils.data.DataLoader(test_set_mix, batch_size =32,
57 shuffle=False, num_workers=8)
58 # get some random training images
59 dataiter = iter(train_loader_mix)
60 images, labels = next(dataiter)
61
62
63 def imshow(img):
64     img = img / 2 + 0.5      # unnormalize
65     npimg = img.numpy()
66     plt.imshow(np.transpose(npimg, (1, 2, 0)))
67 imshow(torchvision.utils.make_grid(images))
68
69
70 model = models.resnet50(weights='ResNet50_Weights.IMGNET1K_V1')
71
72
73
74
75 for param in model.parameters():
76     param.requires_grad = True #False #Set True to train the whole network
77 # Creating final fully connected Layer that accorting to the no of classes we require
78 model.fc = nn.Sequential(nn.Linear(2048, 512),
79                         nn.ReLU(),
80                         nn.Dropout(0.2),
81                         nn.Linear(512, len(class_names))),#,
82                         #nn.LogSoftmax(dim=1))
83
84
85
86 #print(summary(model, input_size=(3, 224, 224)))
87 # Loss and optimizer
88 criterion = nn.CrossEntropyLoss()
89
90 # variable learning rate for different layers and usingg cosine annealing warm
91 # restarts
92 optimizer = optim.SGD([
93     {'params': model.conv1.parameters(), 'lr':1e-4},
94     {'params': model.bn1.parameters(), 'lr': 1e-4},
95     {'params': model.layer1.parameters(), 'lr':1e-4},
96     {'params': model.layer2.parameters(), 'lr':1e-4},
97     {'params': model.layer3.parameters(), 'lr':1e-3},
98     {'params': model.layer4.parameters(), 'lr':1e-3},
99     {'params': model.fc.parameters(), 'lr': 1e-2} # the classifier needs to
100 # learn weights faster
101 ], lr=0.001, weight_decay=0.0005)
102
103
104 if torch.cuda.device_count() > 1:
105     print(f"Using DataParallel Mode with cuda count {torch.cuda.device_count()}")
106     model = torch.nn.DataParallel(model)
107
108
109

```

```

110 model.to(device)
111
112
113 print(summary(model, input_size=(3, 224, 224)))
114
115
116
117 # Restarts the learning rate after every 5 epoch
118 scheduler = lr_scheduler.CosineAnnealingWarmRestarts(
119     optimizer,
120     T_0= 5,
121     T_mult= 1,
122 )
123
124
125 epochs = 200
126 best_acc = 0.0
127 iters = len(train_loader_mix)
128 patience = 500
129 patience_counter = patience
130 best_val_loss = np.inf
131
132
133 train_loss, val_loss = [], []
134
135 import time
136
137
138 print("-----starting training-----")
139 for epoch in range (epochs):
140     start_time = time.time()
141
142     #logger.warning(f"Epoch started: {epoch}")
143     train_loss_epoch = 0.0
144     valid_loss_epoch = 0.0
145     accuracy = 0.0
146
147     # Trainin_the_ model
148     model.train()
149     correct_predictions_train = 0
150     total_samples_train = 0
151
152
153     for i, sample in enumerate (train_loader_mix):
154         inputs, labels = sample
155         # Move data to GPU
156         inputs, labels = inputs.to(device), labels.to(device)
157         # Clear Optimizers
158         optimizer.zero_grad()
159         # Forward Pass
160         logps = model.forward(inputs)
161         # Loss
162         loss = criterion(logps, labels)
163         # Backprop (Calculate Gradients)
164         loss.backward()
165         # Adjust parameters based on gradients
166         optimizer.step()
167         scheduler.step(epoch + i / iters) # if using cosine annealing warm restarts
168         # Add the loss to the running loss
169         train_loss_epoch += loss.item() * inputs.size(0)

```

```

170
171     # Calculate training accuracy
172     _, predicted = torch.max(logps, 1) #_, train_preds = torch.max(logps, 1)
173     correct_predictions_train += torch.sum(predicted == labels).item()#
174     correct_train += torch.sum(train_preds == labels).item()
175     total_samples_train += labels.size(0) #total_train += labels.size(0)
176
177     train_accuracy = correct_predictions_train / total_samples_train# train_accuracy
178     = correct_train / total_train
179
180
181
182     # Update the learning rate scheduler
183     #scheduler.step()
184
185     # Validation
186     model.eval ()
187     correct_predictions = 0
188     total_samples = 0
189     # Tell torch not to calculate gradients
190     with torch.no_grad ():
191         for inputs, labels in val_loader_mix:
192             # Move to device
193             inputs, labels = inputs.to (device), labels.to(device)
194             # Forward pass
195             output = model.forward(inputs)
196             # Calculate Loss
197             val_loss_batch = criterion(output, labels)
198             # Add loss to the validation set's running loss
199             valid_loss_epoch += val_loss_batch.item() * inputs.size(0)
200
201             # Since our model outputs a LogSoftmax, find the real
202             _, predicted = torch.max(output, 1)
203             correct_predictions += torch.sum(predicted == labels).item()
204             total_samples += labels.size(0)
205             accuracy = correct_predictions / total_samples
206
207
208     # Early Stopping
209     if valid_loss_epoch < best_val_loss:
210         best_val_loss = valid_loss_epoch
211         patience_counter = 0
212         torch.save(model.state_dict(),
213                    f'best_model_resnet_modified_300_epochs_28_7_2025_train_accuracy_print_test_run_{test_run}.pth')
214         print(f"Saving the model
215               best_model_resnet_modified_300_epochs_28_7_2025_train_accuracy_print_test_run_{test_run}.pth at epoch {epoch}")
216     else:
217         patience_counter -= 1
218         if patience_counter == 0:
219             print('Early stopping')
220             break
221
222     # Get the average loss for the epoch
223     train_loss_epoch /= len(train_loader_mix.dataset)

```

```

224     valid_loss_epoch /= len(val_loader_mix.dataset)
225
226     # Append the loss and accuracy
227     train_loss.append(train_loss_epoch)
228     val_loss.append(valid_loss_epoch)
229
230
231
232     epoch_time = time.time() - start_time
233     # Print out the information
234     print(f'Epoch {epoch + 1},'
235           f' Training Loss: {train_loss_epoch:.6f},'
236           f' Training Accuracy: {train_accuracy:.6f},'
237           f' Validation Loss: {valid_loss_epoch:.6f},'
238           f' Accuracy: {accuracy:.6f}'
239           f"Time: {epoch_time:.2f} seconds"
240     )
241
242 #logger.warning(f"Training finished")
243
244 print(f"Saving the model
best_model_resnet_modified_300_epochs_28_7_2025_train_accuracy_print_test_run_{test_r
un}_save_after_300_epochs.pth")
245 torch.save(model.state_dict(),
f'best_model_resnet_modified_300_epochs_28_7_2025_train_accuracy_print_test_run_{test_
run}_save_after_300_epochs.pth')
246
247
248 # Plotting the training and validation loss
249 plt.plot(range(1, epochs+1), train_loss, label='Training loss')
250 plt.plot(range(1, epochs+1), val_loss, label='Validation loss')
251 plt.xlabel('Epochs')
252 plt.ylabel('Loss')
253 plt.title(f'Training and Validation Loss_{test_run}')
254 plt.legend()
255 plt.savefig(f"training_loss_curves_test_run_{test_run}.png")
256 plt.show()
257
258 #print(f'Best Validation Accuracy: {best_acc:.6f}')
259
260
261 #####-----Calibration Scores-----#####
262 import numpy as np
263 import torch
264 import torch.nn.functional as F
265 from torch.nn.functional import softmax
266 import numpy
267 from tqdm import tqdm
268
269 alpha = 0.1
270
271
272 #####-----Calculation of
Calibration Scores-----#####
273 def calibration_scores(model, data_loader, device):
274     model.eval()
275

```

```

277     scores_calibration_set = []
278
279     with torch.no_grad():
280         for X, y in tqdm(data_loader, desc="# ##### Proceeding with the
Calculation of Calibration Scores#####"):
281             X = X.to(device)
282             y = y.to(device)
283
284             y = y.cpu().numpy()
285
286             pred_logits = model(X)
287
288             # cal_smx = model(calib_X).softmax(dim=1 ).numpy()
289             cal_smx = softmax(pred_logits, dim=1).cpu().numpy()
290
291             # cal_scores = 1-cal_smx[np . arange(n),cal_labels]
292             scores_per_batch = 1 - cal_smx[np.arange(y.shape[0])], y
293
294             scores_calibration_set.extend(scores_per_batch)
295
296             scores_calibration_set = numpy.array(scores_calibration_set)
297
298     return scores_calibration_set
299 #####-----Calculation of
300 #####-----Prediction Sets-----#####
301
302 #####-----Calculation of
303 def predict_with_conformal_sets(model, data_loader, calibration_scores_array, alpha,
device):
304
305
306     model.eval()
307
308
309     prediction_sets = []
310     true_labels = []
311
312     n = len(calibration_scores_array)
313
314     #get adjusted quantile
315     #q_level = np . ceil((n+1 ) *(1-alpha)) /n
316     q_level = np.ceil((n + 1) * (1 - alpha)) / n
317
318     # qhat = np . quantile(calibration_scores_array, q_level, method='higher')
319     qhat = np.quantile(calibration_scores_array, q_level, method='higher')
320
321     with torch.no_grad():
322         for X, y in tqdm(data_loader, desc="# ##### Proceeding with the
Calculation of Testing Scores#####"):
323
324             X = X.to(device)
325             y = y.cpu().numpy()
326             pred_logits = model(X)
327
328             # cal_smx = model(calib_X).softmax(dim=1 ).numpy()
329             val_smx = softmax(pred_logits, dim=1)
330             val_smx = val_smx.cpu().numpy()
331

```

```

332     # prediction_sets = val_smx >= (1-qhat)
333     masks = val_smx >= (1 - qhat)
334
335     for _ in masks:
336         prediction_sets.append(np.where(_)[0].tolist())
337
338     true_labels.extend(y.tolist())
339
340     return prediction_sets, true_labels
341 #####
342
343
344 #####-----EVALUATION-----#####
345 def evaluate(prediction_sets, labels):
346     number_of_correct_predictions = 0
347     for pred_set, true_label in zip(prediction_sets, labels):
348         if true_label in pred_set:
349             number_of_correct_predictions += 1
350     empirical_test_coverage = number_of_correct_predictions / len(labels)
351     average_prediction_set_size = np.mean([len(s) for s in prediction_sets])
352     return empirical_test_coverage, average_prediction_set_size
353 #####
354
355 # ----- Run Full Conformal Prediction -----
356 scores_calibration = calibration_scores(model=model, data_loader=val_loader_mix,
357 device=device)
357 prediction_sets, true_labels = predict_with_conformal_sets(model=model,
358 data_loader=test_loader_mix, calibration_scores_array=scores_calibration, \
359                                         alpha=alpha,
360                                         device=device)
359 test_set_coverage, avg_size = evaluate(prediction_sets=prediction_sets,
360                                         labels=true_labels)
360
361 import os
362 from datetime import datetime
363 import numpy
364 import pickle
365 import numpy as np
366
367
368 timestamp = datetime.now().strftime("%Y-%m-%d_%H-%M-%S")
369 os.makedirs(f"Results_normal_Conformal_Prediction_experiment_run_{test_run}_{timestamp}",
370 p, exist_ok=True)
370
371 results_directory =
372 f"Results_normal_Conformal_Prediction_experiment_run_{test_run}_{timestamp}"
373
374 import pickle
375 with open(f"{results_directory}/scores_calibration_run_{test_run}_{timestamp}.pkl",
376 "wb") as f:
377     pickle.dump(scores_calibration, f)
378
379 # with open("scores_calibration.pkl", "rb") as f:
380 #     scores_calibration = pickle.load(f)
381

```

```

382
383
384 with
385     open(f"Results_normal_Conformal_Prediction_experiment_run_{test_run}_{timestamp}/logs
386 _experiment_run_{test_run}_{timestamp}.txt", "w") as f:
387         f.write("##### Results\n")
388         f.write(f"Coverage: {test_set_coverage:.3f}\n")
389         f.write(f"Average Prediction Set Size for Test Set: {avg_size:.3f}\n")
390
391
392 print(f"logs file in directory
393 Results_normal_Conformal_Prediction_experiment_run_{test_run}_{timestamp}")
394
395 print("#####Results:#####")
396 print(f"Coverage: {test_set_coverage:.3f}")
397 print(f"Average Prediction Set Size for Test Set: {avg_size:.3f}")
398 print("#####")
399
400 np.save(f"Results_normal_Conformal_Prediction_experiment_run_{test_run}_{timestamp}/c
401 alibration_scores_experiment_run_{test_run}_{timestamp}.npy", scores_calibration)
402
403 with open(f"
404 {results_directory}/prediction_data_experiment_run_{test_run}_{timestamp}.pkl",
405 "wb") as f:
406     pickle.dump({
407         "prediction_sets": prediction_sets,
408         "true_labels": true_labels
409     }, f)
410
411
412 # import numpy as np
413 # import pickle
414
415 # # Load calibration scores
416 # scores = np.load("conformal_results/scores_calibration_2025-08-02_14-30-00.npy")
417
418 # # Load prediction sets and labels
419 # with open("conformal_results/prediction_data_2025-08-02_14-30-00.pkl", "rb") as f:
420 #     data = pickle.load(f)
421 #     prediction_sets = data["prediction_sets"]
422 #     true_labels = data["true_labels"]
423
424
425 #####TO BE
426 CHANGED#####
427 #####
428
429 import numpy as np
430 import matplotlib.pyplot as plt
431 from collections import Counter
432 import pandas as pd
433
434 set_sizes = [len(prediction_set) for prediction_set in prediction_sets]
435 counts = Counter(set_sizes)
436 total = sum(counts.values())
437 sizes = sorted(counts.keys())
438 frequencies = [counts[size] for size in sizes]
439 percentages = [100 * freq / total for freq in frequencies]
440

```

```

432 data = [ {'size': len(prediction_set), 'correct': int(true in prediction_set)} for
433 prediction_set, true in zip(prediction_sets, true_labels)]
434 df = pd.DataFrame(data)
435 coverage_stats = df.groupby('size').agg({'correct': ['mean', 'sum', 'count']})
436 coverage_stats.columns = ['coverage', 'correct_count', 'total_count']
437 coverage_stats = coverage_stats.reindex(sizes, fill_value=0)
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}%', xy=(bar.get_x() + bar.get_width() / 2, height),
459                 xytext=(0, 8),
460                 textcoords="offset points",
461                 ha='center', va='bottom', fontsize=9, color='darkorange')
462
463 mean_size = np.mean(set_sizes)
464 median_size = np.median(set_sizes)
465 ax.axvline(mean_size, color='gray', linestyle='--', label=f'Mean = {mean_size:.2f}')
466 ax.axvline(median_size, color='orange', linestyle='--', label=f'Median = {median_size:.2f}')
467
468 ax.set_xlabel("Prediction Set Size")
469 ax.set_ylabel("Frequency")
470 ax.set_title("Prediction Set Size Distribution with Coverage (Test Set)")
471 ax.set_xticks(sizes)
472 ax.grid(axis='y', linestyle='--', alpha=0.6)
473 ax.legend(loc='upper right')
474
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"))
479
480
481 plt.show()
482
483
484

```

```

485 #####script#####
486 from collections import Counter
487 import matplotlib.pyplot as plt
488 from collections import defaultdict
489 import numpy as np
490 from scipy.stats import pearsonr, spearmanr
491
492
493 def feature_stratified_coverage_by_class(prediction_sets, true_labels, alpha=0.1,
494     class_names=None):
495
496     class_to_indices = defaultdict(list)
497     for idx, label in enumerate(true_labels):
498         class_to_indices[label].append(idx)
499
500     class_coverages = {}
501     for cls, indices in class_to_indices.items():
502         covered = sum(true_labels[i] in prediction_sets[i] for i in indices)
503         coverage = covered / len(indices)
504         class_coverages[cls] = coverage
505
506     print("\n Feature-Stratified Coverage (by True Class Label):")
507     for cls, cov in sorted(class_coverages.items()):
508         name = class_names[cls] if class_names else str(cls)
509         print(f" Class {name:<20}: Coverage = {cov:.3f} (Target ≥ {1 - alpha:.2f})")
510
511     fsc_metric = min(class_coverages.values())
512     print(f"\n FSC Metric (minimum class-wise coverage): {fsc_metric:.3f}")
513
514     return fsc_metric, class_coverages
515
516
517 fsc_metric, per_class_cov = feature_stratified_coverage_by_class(
518     prediction_sets,
519     true_labels,
520     alpha=alpha,
521     class_names=class_names
522 )
523
524
525 def plot_class_wise_coverage(class_coverages, alpha=0.1, class_names=None,
526     title="Class-wise Coverage (FSC)"):
527
528     classes = list(class_coverages.keys())
529     coverages = [class_coverages[c] for c in classes]
530
531     # Get display names
532     labels = [class_names[c] if class_names else str(c) for c in classes]
533
534     # Plot
535     plt.figure(figsize=(12, 6))
536     bars = plt.bar(labels, coverages, color="pink", edgecolor="black")
537
538     # Target line
539     target = 1 - alpha
540     plt.axhline(target, color='red', linestyle='--', label=f"Target Coverage
({target:.2f})")

```

```

541 # Highlight underperforming classes
542 for bar, cov in zip(bars, coverages):
543     if cov < target:
544         bar.set_color('salmon')
545
546 plt.xticks(rotation=45, ha="right")
547 plt.ylim(0, 1.05)
548 plt.ylabel("Coverage")
549 plt.title(title)
550 plt.legend()
551 plt.tight_layout()
552 filename = f"Class_wise_Coverage_(FSC)_{test_run}_{timestamp}.png"
553 plt.savefig(os.path.join(results_directory, filename))
554 plt.show()
555
556 fsc_metric, class_coverages = feature_stratified_coverage_by_class(
557     prediction_sets,
558     true_labels,
559     alpha=alpha,
560     class_names=class_names
561 )
562
563 plot_class_wise_coverage(class_coverages, alpha=alpha, class_names=class_names)
564
565
566
567
568
569 def plot_coverage_vs_class_frequency(class_coverages, true_labels, alpha=0.1,
570                                         class_names=None, title="Coverage vs Class Frequency"):
571
572
573     class_indices = sorted(class_coverages.keys())
574     coverages = [class_coverages[c] for c in class_indices]
575
576
577     total = len(true_labels)
578     counts = Counter(true_labels)
579     frequencies = [counts[c] / total for c in class_indices]
580
581
582     labels = [class_names[c] if class_names else str(c) for c in class_indices]
583     x = np.arange(len(class_indices))
584
585     fig, ax1 = plt.subplots(figsize=(12, 6))
586
587
588     bars1 = ax1.bar(x - 0.2, coverages, width=0.4, label="Coverage", color='pink',
589                     edgecolor='black')
590     ax1.axhline(1 - alpha, color='red', linestyle='--', label=f"Target Coverage ({1 -
591     alpha:.2f})")
592     ax1.set_ylim(0, 1.05)
593     ax1.set_ylabel("Coverage", color='pink')
594     ax1.tick_params(axis='y', labelcolor='pink')
595
596     # Highlight under-covered bars
597     for bar, cov in zip(bars1, coverages):
598         if cov < (1 - alpha):
599             bar.set_color('salmon')

```

```

598
599 # Twin y-axis for frequencies
600 ax2 = ax1.twinx()
601 bars2 = ax2.bar(x + 0.2, frequencies, width=0.4, label="Class Frequency",
602 color='gray', alpha=0.5)
603 ax2.set_ylabel("Class Frequency", color='gray')
604 ax2.tick_params(axis='y', labelcolor='gray')
605
606 plt.xticks(x, labels, rotation=45, ha="right")
607
608 # Title and legend
609 plt.title(title)
610 fig.legend(loc='upper right', bbox_to_anchor=(0.85, 0.85))
611 plt.tight_layout()
612 filename = f"plot_coverage_vs_class_frequency_{test_run}_{timestamp}.png"
613 plt.savefig(os.path.join(results_directory, filename))
614 plt.show()
615
616
617 # Step 1: Compute coverage by class
618 fsc_metric, class_coverages = feature_stratified_coverage_by_class(
619     prediction_sets,
620     true_labels,
621     alpha=alpha,
622     class_names=class_names
623 )
624
625 # Step 2: Compare against class frequency
626 plot_coverage_vs_class_frequency(class_coverages, true_labels, alpha=alpha,
627                                     class_names=class_names)
628
629 def plot_coverage_vs_class_frequency_with_correlation(class_coverages, true_labels,
630 alpha=0.1, class_names=None, title="Coverage vs Class Frequency"):
631
632
633     class_indices = sorted(class_coverages.keys())
634     coverages = [class_coverages[c] for c in class_indices]
635
636     # Compute frequencies
637     total = len(true_labels)
638     counts = Counter(true_labels)
639     frequencies = [counts[c] / total for c in class_indices]
640
641     # Class labels
642     labels = [class_names[c] if class_names else str(c) for c in class_indices]
643     x = np.arange(len(class_indices))
644
645     # Correlation analysis
646     pearson_corr, pearson_p = pearsonr(frequencies, coverages)
647     spearman_corr, spearman_p = spearmanr(frequencies, coverages)
648
649     print(f"\n Correlation between class frequency and coverage:")
650     print(f" - Pearson r = {pearson_corr:.3f}, p = {pearson_p:.3e}")
651     print(f" - Spearman p = {spearman_corr:.3f}, p = {spearman_p:.3e}")
652
653     # Plotting
654     fig, ax1 = plt.subplots(figsize=(12, 6))

```

```

655
656     bars1 = ax1.bar(x - 0.2, coverages, width=0.4, label="Coverage", color='pink',
657     edgecolor='black')
658     ax1.axhline(1 - alpha, color='red', linestyle='--', label=f"Target Coverage ({1 -
659     alpha:.2f})")
660     ax1.set_ylim(0, 1.05)
661     ax1.set_ylabel("Coverage", color='pink')
662     ax1.tick_params(axis='y', labelcolor='pink')
663
664     # Highlight under-covered classes
665     for bar, cov in zip(bars1, coverages):
666         if cov < (1 - alpha):
667             bar.set_color('salmon')
668
669     # Frequency bars on second axis
670     ax2 = ax1.twinx()
671     bars2 = ax2.bar(x + 0.2, frequencies, width=0.4, label="Class Frequency",
672     color='gray', alpha=0.5)
673     ax2.set_ylabel("Class Frequency", color='gray')
674     ax2.tick_params(axis='y', labelcolor='gray')
675
676     # X-axis and titles
677     plt.xticks(x, labels, rotation=45, ha="right")
678     plt.title(title)
679     fig.legend(loc='upper right', bbox_to_anchor=(0.85, 0.85))
680     plt.tight_layout()
681     filename =
682         f"plot_coverage_vs_class_frequency_with_correlation_{test_run}_{timestamp}.png"
683     plt.savefig(os.path.join(results_directory, filename))
684     plt.show()
685
686
687
688 )
689
690
691
692 def size_stratified_coverage(prediction_sets, true_labels, alpha=0.1, bins=[[1], [2],
693 [3, 4, 5], list(range(6, 100))]):
694     """
695         Computes size-stratified coverage (SSC): coverage within prediction set size
696         bins.
697
698     Args:
699         prediction_sets (List[List[int]]): output sets from conformal_predict
700         true_labels (List[int]): ground-truth labels
701         alpha (float): target error level (e.g., 0.1)
702         bins (List[List[int]]): list of prediction set size groups
703
704     Returns:
705         ssc (float): min group coverage (worst-case)
706         bin_coverages (dict): {bin_range: coverage}
707     """
708     bin_groups = defaultdict(list)
709
710     # Assign examples to bins based on set size

```

```

709     for i, pred_set in enumerate(prediction_sets):
710         set_size = len(pred_set)
711         for bin_range in bins:
712             if set_size in bin_range:
713                 bin_groups[str(bin_range)].append(i)
714                 break
715
716     bin_coverages = {}
717     for bin_key, indices in bin_groups.items():
718         covered = sum(true_labels[i] in prediction_sets[i] for i in indices)
719         coverage = covered / len(indices) if indices else 0
720         bin_coverages[bin_key] = coverage
721
722     print("\n Size-Stratified Coverage (SSC):")
723     for bin_key, cov in bin_coverages.items():
724         print(f" Set Size {bin_key}:10s: Coverage = {cov:.3f} (Target ≥ {1 - alpha:.2f})")
725
726     ssc = min(bin_coverages.values())
727     print(f"\n SSC Metric (min bin-wise coverage): {ssc:.3f}")
728     return ssc, bin_coverages
729
730
731 ssc_metric, bin_coverages = size_stratified_coverage(
732     prediction_sets,
733     true_labels,
734     alpha=alpha,
735     bins=[[1], [2], [3], list(range(4, 100))])
736 )
737
738
739
740 def plot_ssc_coverage(bin_coverages, alpha=0.1, title="Size-Stratified Coverage"):
741     """
742     Plots coverage per prediction set size bin.
743
744     Args:
745         bin_coverages (dict): mapping from bin label (e.g. '[1]', '[2, 3]') to
746         coverage
747         alpha (float): target error level
748         title (str): plot title
749     """
750     bin_labels = list(bin_coverages.keys())
751     coverages = [bin_coverages[k] for k in bin_labels]
752
753     plt.figure(figsize=(10, 5))
754     bars = plt.bar(bin_labels, coverages, color='pink', edgecolor='black')
755
756     # Target line
757     target = 1 - alpha
758     plt.axhline(target, color='red', linestyle='--', label=f"Target Coverage
({target:.2f})")
759
760     # Highlight low-coverage bins
761     for bar, cov in zip(bars, coverages):
762         if cov < target:
763             bar.set_color('grey')
764
765     plt.xticks(rotation=45, ha='right')
766     plt.ylim(0, 1.05)

```

```

766     plt.ylabel("Coverage")
767     plt.title(title)
768     plt.legend()
769     plt.tight_layout()
770     filename = f"plot_ssc_coverage_{test_run}_{timestamp}.png"
771     plt.savefig(os.path.join(results_directory, filename))
772     plt.show()
773
774
775 ssc_metric, bin_coverages = size_stratified_coverage(
776     prediction_sets,
777     true_labels,
778     alpha=alpha,
779     bins=[[1], [2], [3], list(range(4, 100))])
780 )
781
782 plot_ssc_coverage(bin_coverages, alpha=alpha)
783
784
785
786 def plot_prediction_set_size_by_class(prediction_sets, true_labels, class_names=None,
787                                         results_dir="", test_run="", timestamp ""):
788
789     # Group prediction set sizes by true class
790     size_per_class = defaultdict(list)
791     for pred_set, true_label in zip(prediction_sets, true_labels):
792         size_per_class[true_label].append(len(pred_set))
793
794     # Sort classes
795     sorted_classes = sorted(size_per_class.keys())
796     means = [np.mean(size_per_class[c]) for c in sorted_classes]
797     stds = [np.std(size_per_class[c]) for c in sorted_classes]
798     labels = [class_names[c] if class_names else str(c) for c in sorted_classes]
799
800     # Plot
801     plt.figure(figsize=(12, 6))
802     plt.bar(labels, means, yerr=stds, capsized=5, color="lightgreen",
803             edgecolor="black")
804     plt.ylabel("Avg Prediction Set Size")
805     plt.xlabel("Class")
806     plt.title("Average Prediction Set Size by Class")
807     plt.xticks(rotation=45, ha="right")
808     plt.grid(axis="y", linestyle="--", alpha=0.5)
809     plt.tight_layout()
810
811     # Save
812     filename = f"prediction_set_size_by_class_{test_run}_{timestamp}.png"
813     plt.savefig(os.path.join(results_directory, filename))
814     plt.show()
815
816     plot_prediction_set_size_by_class(
817         prediction_sets=prediction_sets,
818         true_labels=true_labels,
819         class_names=class_names,
820         results_dir=results_directory,
821         test_run=test_run,
822         timestamp=timestamp
823     )

```

```

824
825
826
827
828 def plot_calibration_score_distribution(calibration_scores, results_dir="", test_run="", timestamp ""):
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
838     filename = f"calibration_score_distribution_{test_run}_{timestamp}.png"
839     plt.savefig(os.path.join(results_directory, filename))
840
841     plt.show()
842
843
844 plot_calibration_score_distribution(
845     calibration_scores=scores_calibration,
846     results_dir=results_directory,
847     test_run=test_run,
848     timestamp=timestamp
849 )
850
851
852 import json
853
854 metrics = {
855     "test_run": test_run,
856     "timestamp": timestamp,
857     "alpha": alpha,
858     "test_set_coverage": test_set_coverage,
859     "average_prediction_set_size": avg_size,
860     "FSC": fsc_metric,
861     "SSC": ssc_metric,
862     "class_coverages": {str(k): float(v) for k, v in class_coverages.items()},
863     "bin_coverages": bin_coverages
864 }
865
866 with open(os.path.join(results_directory,
867     f"metrics_summary_{test_run}_{timestamp}.json"), "w") as f:
868     json.dump(metrics, f, indent=4)
869
870
871
872 ######
873 ######
874 #####

```