```python
1  """
2     evaluation_pipeline_steps: Implementations of p
   ipeline steps to calculate eval metrics and print t
   ables.
3  """
4
5  import logging
6
7  import numpy as np
8  import preprocessing
9  from pipeline_framework import PipelineState, Pipel
   ineStep
10
11
12 class EvaluatePredictionsPipelineStep(PipelineSte
   p):
13     step_key = "evaluate_predictions"
14
15     def __init__(
16         self, config: dict, state: PipelineState, l
   ogger: logging.Logger, split: str
17     ):
18         super().__init__(config, state, logger)
19
20         self.split = split
21         self.label_strategy = preprocessing.LabelSt
   rategy[self.config["label_strategy"]]
22
23         self.evaluation_yaml_path = (
24             self.state.run_base_dir / f"evaluation_
   {self.split}.yaml"
25         )
26
27     def input_ready(self, input: dict) -> bool:
28         return "examples" in input.keys() and any(
29             ["Prediction" in example.keys() for exa
   mple in input["examples"]]
30         )
31
32     def execute(self, input: dict, output: dict) ->
   None:
33         # Collect relevant data
34         split_examples = [ex for ex in input["examp
   les"] if ex["Split"] == self.split]
35
36         # Overall evaluation
37         overall_result = self._evaluate_subset(spli
   t_examples, input["label_encoding"])
38
39         # Per CWE
40         weakness_ids = [cwe["WeaknessID"] for cwe i
   n input["cwes"]]
41         per_cwe_result = self._evaluate_subsets(
42             split_examples,
43             input["label_encoding"],
44             set(weakness_ids),
45             lambda ex, subset: ex["Testcase"]["Weak
   ness"]["WeaknessID"] == subset,
46         )
47
```

```python
1  """
2     evaluation_pipeline_steps: Implementations of p
   ipeline steps to calculate eval metrics and print t
   ables.
3  """
4
5  import logging
6
7  import numpy as np
8  import preprocessing
9  from pipeline_framework import PipelineState, Pipel
   ineStep
10
11
12 class EvaluatePredictionsPipelineStep(PipelineSte
   p):
13     step_key = "evaluate_predictions"
14
15     def __init__(
16         self, config: dict, state: PipelineState, l
   ogger: logging.Logger, split: str
17     ):
18         super().__init__(config, state, logger)
19
20         self.split = split
21         self.label_strategy = preprocessing.LabelSt
   rategy[self.config["label_strategy"]]
22
23         self.evaluation_yaml_path = (
24             self.state.run_base_dir / f"evaluation_
   {self.split}.yaml"
25         )
26
27     def input_ready(self, input: dict) -> bool:
28         return "examples" in input.keys() and any(
29             ["Prediction" in example.keys() for exa
   mple in input["examples"]]
30         )
31
32     def execute(self, input: dict, output: dict) ->
   None:
33         # Collect relevant data
34         split_examples = [ex for ex in input["examp
   les"] if ex["Split"] == self.split]
35
36         # Overall evaluation
37         overall_result = self._evaluate_subset(spli
   t_examples, input["label_encoding"])
38
39         if not self.state.malware:
40             # Per CWE
41             weakness_ids = [cwe["WeaknessID"] for c
   we in input["cwes"]]
42             per_cwe_result = self._evaluate_subsets
   (
43                 split_examples,
44                 input["label_encoding"],
45                 set(weakness_ids),
46                 lambda ex, subset: ex["Testcase"]
   ["Weakness"]["WeaknessID"] == subset,
47             )
48
```

Left side:

```python
48          # Per Flow Variant
49          flow_variants = set(
50              [testcase["FlowVariant"] for testcase i
n input["testcases"]]
51          )
52          per_flow_variant_result = self._evaluate_su
bsets(
53              split_examples,
54              input["label_encoding"],
55              flow_variants,
56              lambda ex, subset: ex["Testcase"]["Flow
Variant"] == subset,
57          )
58
59          bvdetector_subsets = {
60              "MC": {
61                  120,
62                  124,
63                  126,
64                  127,
65                  129,
66                  134,
67                  170,
68                  415,
69                  416,
70                  590,
71                  761,
72                  785,
73                  805,
74                  806,
75                  822,
76                  824,
77                  843,
78              },
79              "NH": {190, 191, 194, 195, 196, 197, 36
9, 682, 839},
80          }
81          bvdetector_subsets["MC&NH"] = bvdetector_su
bsets["MC"].union(
82              bvdetector_subsets["NH"]
83          )
84
85          bvdetector_subset_result = self._evaluate_s
ubsets(
86              split_examples,
87              input["label_encoding"],
88              bvdetector_subsets.keys(),
89              lambda ex, subset: ex["Testcase"]["Weak
ness"]["WeaknessID"]
90              in bvdetector_subsets[subset],
91          )
92
93          evaluation = {
94              "Overall": overall_result,
95              "PerCWE": per_cwe_result,
96              "PerFlowVariant": per_flow_variant_resu
lt,
97              "BVDetector": bvdetector_subset_result,
98          }
```

Right side:

```python
49          # Per Flow Variant
50          flow_variants = set(
51              [testcase["FlowVariant"] for testca
se in input["testcases"]]
52          )
53          per_flow_variant_result = self._evaluat
e_subsets(
54              split_examples,
55              input["label_encoding"],
56              flow_variants,
57              lambda ex, subset: ex["Testcase"]
["FlowVariant"] == subset,
58          )
59
60          bvdetector_subsets = {
61              "MC": {
62                  120,
63                  124,
64                  126,
65                  127,
66                  129,
67                  134,
68                  170,
69                  415,
70                  416,
71                  590,
72                  761,
73                  785,
74                  805,
75                  806,
76                  822,
77                  824,
78                  843,
79              },
80              "NH": {190, 191, 194, 195, 196, 19
7, 369, 682, 839},
81          }
82          bvdetector_subsets["MC&NH"] = bvdetecto
r_subsets["MC"].union(
83              bvdetector_subsets["NH"]
84          )
85
86          bvdetector_subset_result = self._evalua
te_subsets(
87              split_examples,
88              input["label_encoding"],
89              bvdetector_subsets.keys(),
90              lambda ex, subset: ex["Testcase"]
["Weakness"]["WeaknessID"]
91              in bvdetector_subsets[subset],
92          )
93
94          evaluation = {
95              "Overall": overall_result,
96              "PerCWE": per_cwe_result,
97              "PerFlowVariant": per_flow_variant_
result,
98              "BVDetector": bvdetector_subset_res
ult,
99          }
100     else:
101          evaluation = {
102              "Overall": overall_result,
103              "PerCWE": overall_result,
104              "PerFlowVariant": overall_result,
105              "BVDetector": overall_result,
```

```
              }

          self.logger.info(f"Evaluation: {evaluatio
n}")

          # Save to YAML
          with open(self.evaluation_yaml_path, "w") a
s evaluation_yaml:
              import yaml

              yaml.dump(evaluation, evaluation_yaml)

          if "evaluation" in output.keys():
              output["evaluation"][self.split] = eval
uation
          else:
              output["evaluation"] = {self.split: eva
luation}

      def _evaluate_subsets(
          self,
          examples: list[dict],
          label_encoding: dict,
          subsets: set[str],
          subset_filter,
      ):
          subsets_results = {}

          for subset in subsets:
              subset_examples = [ex for ex in example
s if subset_filter(ex, subset)]
              subsets_results[subset] = self._evaluat
e_subset(
                  subset_examples, label_encoding
              )

          return subsets_results

      def _evaluate_subset(self, examples: list[dic
t], label_encoding: dict):
          if self.label_strategy == preprocessing.Lab
elStrategy.BINARYCLASSIFICATION:
              return self._evaluate_subset_binary(exa
mples, label_encoding)
          else:
              return self._evaluate_subset_multiclass
(examples, label_encoding)

      def _evaluate_subset_binary(self, examples: lis
t[dict], label_encoding: dict):
          confusion_matrix = self._calculate_confusio
n_matrix(examples, label_encoding)

          true_positives = float(confusion_matrix[1,
 1])
          false_positives = float(confusion_matrix[0,
1])
          false_negatives = float(confusion_matrix[1,
0])
          true_negatives = float(confusion_matrix[0,
 0])
          total_examples = (
              true_positives + false_positives + fals
e_negatives + true_negatives
          )

          # Accuracy
          if total_examples > 0:
```

```
149            accuracy = (true_positives + true_negat
        ives) / (
150                true_positives + false_positives +
         false_negatives + true_negatives
151            )
152        else:
153            accuracy = None
154
155        # Precision
156        if (true_positives + false_positives) > 0:
157            precision = true_positives / (true_posi
        tives + false_positives)
158        else:
159            precision = None
160
161        # Recall, TPR, FNR
162        if (true_positives + false_negatives) > 0:
163            recall = true_positives / (true_positiv
        es + false_negatives)
164            true_positive_rate = true_positives /
         (true_positives + false_negatives)
165            false_negative_rate = false_negatives /
         (false_negatives + true_positives)
166        else:
167            recall = None
168            true_positive_rate = None
169            false_negative_rate = None
170
171        # FPR, TNR
172        if (true_negatives + false_positives) > 0:
173            false_positive_rate = false_positives /
         (false_positives + true_negatives)
174            true_negative_rate = true_negatives /
         (true_negatives + false_positives)
175        else:
176            false_positive_rate = None
177            true_negative_rate = None
178
179        # F1
180        if precision is not None and recall is not
         None:
181            if (precision + recall) > 0:
182                f1 = 2 * precision * recall / (prec
        ision + recall)
183            else:
184                f1 = 0.0
185        else:
186            f1 = None
187
188        return {
189            "Accuracy": accuracy,
190            "Precision": precision,
191            "Recall": recall,
192            "F1": f1,
193            "TotalExamples": int(total_examples),
194            "TotalPositives": int(true_positives +
         false_negatives),
195            "TotalNegatives": int(true_negatives +
         false_positives),
196            "TruePositiveRate": true_positive_rate,
197            "FalsePositiveRate": false_positive_rat
        e,
198            "FalseNegativeRate": false_negative_rat
        e,
199            "TrueNegativeRate": true_negative_rate,
200        }
```

```python
201
202     def _evaluate_subset_multiclass(self, examples:
    list[dict], label_encoding: dict):
203         confusion_matrix = self._calculate_confusio
    n_matrix(examples, label_encoding)
204
205         total_examples = float(confusion_matrix.sum
    ())
206         if total_examples > 0:
207             accuracy = float(confusion_matrix.diago
    nal().sum()) / total_examples
208         else:
209             return None
210
211         return {"Accuracy": accuracy, "TotalExample
    s": total_examples}
212
213     def _calculate_confusion_matrix(self, examples:
    list[dict], label_encoding: dict):
214         label_count = len(label_encoding.keys())
215         confusion_matrix = np.zeros((label_count, l
    abel_count), dtype=np.int64)
216
217         for example in examples:
218             prediction = example["Prediction"]
219             label = example["Label"]
220
221             label_idx = label_encoding[label]
222             prediction_idx = label_encoding[predict
    ion]
223
224             confusion_matrix[label_idx, prediction_
    idx] += 1
225
226         return confusion_matrix
227
228
229 class PrintEvalulationTablePipelineStep(PipelineSte
    p):
230     step_key = "print_evaluation_table"
231
232     def __init__(
233         self,
234         config: dict,
235         state: PipelineState,
236         logger: logging.Logger,
237         collection: str,
238         split: str,
239         keys: list[str],
240     ):
241         super().__init__(config, state, logger)
242
243         self.collection = collection
244         self.keys = keys
245         self.split = split
246
247     def input_ready(self, input: dict) -> bool:
248         return (
249             "evaluation" in input.keys()
250             and self.split in input["evaluation"].k
    eys()
251             and self.collection in input["evaluatio
    n"][self.split].keys()
252         )
253
254     def execute(self, input: dict, output: dict) ->
    None:
```

```python
209
210     def _evaluate_subset_multiclass(self, examples:
    list[dict], label_encoding: dict):
211         confusion_matrix = self._calculate_confusio
    n_matrix(examples, label_encoding)
212
213         total_examples = float(confusion_matrix.sum
    ())
214         if total_examples > 0:
215             accuracy = float(confusion_matrix.diago
    nal().sum()) / total_examples
216         else:
217             return None
218
219         return {"Accuracy": accuracy, "TotalExample
    s": total_examples}
220
221     def _calculate_confusion_matrix(self, examples:
    list[dict], label_encoding: dict):
222         label_count = len(label_encoding.keys())
223         confusion_matrix = np.zeros((label_count, l
    abel_count), dtype=np.int64)
224
225         for example in examples:
226             prediction = example["Prediction"]
227             label = example["Label"]
228
229             label_idx = label_encoding[label]
230             prediction_idx = label_encoding[predict
    ion]
231
232             confusion_matrix[label_idx, prediction_
    idx] += 1
233
234         return confusion_matrix
235
236
237 class PrintEvalulationTablePipelineStep(PipelineSte
    p):
238     step_key = "print_evaluation_table"
239
240     def __init__(
241         self,
242         config: dict,
243         state: PipelineState,
244         logger: logging.Logger,
245         collection: str,
246         split: str,
247         keys: list[str],
248     ):
249         super().__init__(config, state, logger)
250
251         self.collection = collection
252         self.keys = keys
253         self.split = split
254
255     def input_ready(self, input: dict) -> bool:
256         return (
257             "evaluation" in input.keys()
258             and self.split in input["evaluation"].k
    eys()
259             and self.collection in input["evaluatio
    n"][self.split].keys()
260         )
261
262     def execute(self, input: dict, output: dict) ->
    None:
```

```
255        evaluation_of_collection = input["evaluatio
   n"][self.split][self.collection]
256        subsets = list(sorted(evaluation_of_collect
   ion.keys()))
257
258        table_tuples = []
259
260        # Compose table header
261        table_tuples += [("Subset",) + tuple(self.k
   eys)]
262        table_tuples += [tuple(["-"] * (len(self.ke
   ys) + 1))]
263
264        # Compose table rows
265        for subset in subsets:
266            row = [str(subset)]
267            for key in self.keys:
268                if isinstance(evaluation_of_collect
   ion[subset][key], float):
269                    row += (f"{evaluation_of_collec
   tion[subset][key]:0.4f}",)
270                else:
271                    row += (str(evaluation_of_colle
   ction[subset][key]),)
272
273            table_tuples.append(row)
274
275        # Convert table to strings
276        table_strings = []
277        for table_tuple in table_tuples:
278            table_strings.append(
279                "| " + " | ".join([f"{str(el):15s}"
   for el in table_tuple]) + " |"
280            )
281
282        # Print table strings to logger
283        self.logger.info(
284            f"Evaluation table of split {self.spli
   t} (collection {self.collection}):\n"
285            + "\n".join(table_strings)
286        )
287
288
289 __all__ = ["EvaluatePredictionsPipelineStep", "Prin
   tEvalulationTablePipelineStep"]
290
```