

```

1 # coding=utf-8
2 # Copyright 2020-present the HuggingFace Inc. team.
3 #
4 # Licensed under the Apache License, Version 2.0 (the "License");
5 # you may not use this file except in compliance with the License.
6 # You may obtain a copy of the License at
7 #
8 #     http://www.apache.org/licenses/LICENSE-2.0
9 #
10 # Unless required by applicable law or agreed to in writing, software
11 # distributed under the License is distributed on an "AS IS" BASIS,
12 # WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
13 # See the License for the specific language governing permissions and
14 # limitations under the License.
15 """
16 Torch utilities for the Trainer class.
17 """
18
19 import datetime
20 import json
21 import math
22 import os
23 import warnings
24 from contextlib import contextmanager
25 from dataclasses import dataclass
26 from typing import Dict, Iterator, List, Optional, Union
27
28 import numpy as np
29 import torch
30 from packaging import version
31 from torch.utils.data.dataset import Dataset, IterableDataset
32 from torch.utils.data.distributed import DistributedSampler
33 from torch.utils.data.sampler import RandomSampler, Sampler
34
35 from .file_utils import is_sagemaker_dp_enabled, is_sagemaker_mp_enabled, is_torch_tpu_available
36 from .utils import logging
37
38 if is_sagemaker_dp_enabled():
39     import smdistributed.dataparallel.torch.distributed as dist
40 else:
41     import torch.distributed as dist
42
43 if is_torch_tpu_available():
44     import torch_xla.core.xla_model as xm
45
46 # this is used to suppress an undesired warning emitted by pytorch versions 1.4.2-1.7.0
47 try:
48     from torch.optim.lr_scheduler import SAVE_STATE_WARNING
49 except ImportError:
50     SAVE_STATE_WARNING = ""
51
52 logger = logging.get_logger(__name__)
53
54
55 def torch_pad_and_concatenate(tensor1, tensor2, padding_index=-100):
56     """Concatenates `tensor1` and `tensor2` on first axis, applying padding on the second if necessary."""
57     if len(tensor1.shape) == 1 or tensor1.shape[1] == tensor2.shape[1]:
58         return torch.cat((tensor1, tensor2), dim=0)
59
60     # Let's figure out the new shape
61     new_shape = (tensor1.shape[0] + tensor2.shape[0], max(tensor1.shape[1], tensor2.shape[1])) + tensor1.shape[2:]
62
63     # Now let's fill the result tensor
64     result = tensor1.new_full(new_shape, padding_index)
65     result[: tensor1.shape[0], : tensor1.shape[1]] = tensor1
66     result[tensor1.shape[0] :, : tensor2.shape[1]] = tensor2
67     return result
68
69
70 def numpy_pad_and_concatenate(array1, array2, padding_index=-100):
71     """Concatenates `array1` and `array2` on first axis, applying padding on the second if necessary."""
72     if len(array1.shape) == 1 or array1.shape[1] == array2.shape[1]:
73         return np.concatenate((array1, array2), axis=0)
74
75     # Let's figure out the new shape
76     new_shape = (array1.shape[0] + array2.shape[0], max(array1.shape[1], array2.shape[1])) + array1.shape[2:]
77
78     # Now let's fill the result tensor
79     result = np.full_like(array1, padding_index, shape=new_shape)
80     result[: array1.shape[0], : array1.shape[1]] = array1
81     result[array1.shape[0] :, : array2.shape[1]] = array2
82     return result
83
84
85 def nested_concat(tensors, new_tensors, padding_index=-100):
86     """
87     Concat the `new_tensors` to `tensors` on the first dim and pad them on the second if needed. Works for tensors or
88     nested list/tuples of tensors.
89     """
90     assert type(tensors) == type(
91         new_tensors
92     ), f"Expected `tensors` and `new_tensors` to have the same type but found {type(tensors)} and {type(new_tensors)}."
93     if isinstance(tensors, (list, tuple)):
94         return type(tensors)(nested_concat(t, n, padding_index=padding_index) for t, n in zip(tensors, new_tensors))
95     elif isinstance(tensors, torch.Tensor):
96         return torch_pad_and_concatenate(tensors, new_tensors, padding_index=padding_index)
97     elif isinstance(tensors, np.ndarray):
98         return numpy_pad_and_concatenate(tensors, new_tensors, padding_index=padding_index)
99     else:
100         raise TypeError(f"Unsupported type for concatenation: got {type(tensors)}")
101
102
103 def find_batch_size(tensors):
104     """
105     Find the first dimension of a tensor in a nested list/tuple/dict of tensors.
106     """
107     if isinstance(tensors, (list, tuple)):
108         for t in tensors:
109             result = find_batch_size(t)
110             if result is not None:
111                 return result
112     elif isinstance(tensors, dict):
113         for key, value in tensors.items():
114             result = find_batch_size(value)
115             if result is not None:
116                 return result
117     elif isinstance(tensors, torch.Tensor):

```

```

120         return tensors.shape[0] if len(tensors.shape) >= 1 else None
121     elif isinstance(tensors, np.ndarray):
122         return tensors.shape[0] if len(tensors.shape) >= 1 else None
123
124
125 def nested_numpyify(tensors):
126     "Numpify `tensors` (even if it's a nested list/tuple of tensors)."
127     if isinstance(tensors, (list, tuple)):
128         return type(tensors)(nested_numpyify(t) for t in tensors)
129     return tensors.cpu().numpy()
130
131
132 def nested_detach(tensors):
133     "Detach `tensors` (even if it's a nested list/tuple of tensors)."
134     if isinstance(tensors, (list, tuple)):
135         return type(tensors)(nested_detach(t) for t in tensors)
136     return tensors.detach()
137
138
139 def nested_xla_mesh_reduce(tensors, name):
140     if is_torch_tpu_available():
141         import torch_xla.core.xla_model as xm
142
143         if isinstance(tensors, (list, tuple)):
144             return type(tensors)(nested_xla_mesh_reduce(t, f"{name}_{i}") for i, t in enumerate(tensors))
145         return xm.mesh_reduce(name, tensors, torch.cat)
146     else:
147         raise ImportError("Torch xla must be installed to use `nested_xla_mesh_reduce`")
148
149
150 def distributed_concat(tensor: "torch.Tensor", num_total_examples: Optional[int] = None) -> torch.Tensor:
151     try:
152         if isinstance(tensor, (tuple, list)):
153             return type(tensor)(distributed_concat(t, num_total_examples) for t in tensor)
154         output_tensors = [tensor.clone() for _ in range(dist.get_world_size())]
155         dist.all_gather(output_tensors, tensor)
156         concat = torch.cat(output_tensors, dim=0)
157
158         # truncate the dummy elements added by SequentialDistributedSampler
159         if num_total_examples is not None:
160             concat = concat[:num_total_examples]
161         return concat
162     except AssertionError:
163         raise AssertionError("Not currently using distributed training")
164
165
166 def distributed_broadcast_scalars(
167     scalars: List[Union[int, float]], num_total_examples: Optional[int] = None
168 ) -> torch.Tensor:
169     try:
170         tensorized_scalar = torch.tensor(scalars).cuda()
171         output_tensors = [tensorized_scalar.clone() for _ in range(dist.get_world_size())]
172         dist.all_gather(output_tensors, tensorized_scalar)
173         concat = torch.cat(output_tensors, dim=0)
174
175         # truncate the dummy elements added by SequentialDistributedSampler
176         if num_total_examples is not None:
177             concat = concat[:num_total_examples]
178         return concat
179     except AssertionError:
180         raise AssertionError("Not currently using distributed training")
181
182
183 def reissue_pt_warnings(caught_warnings):
184     # Reissue warnings that are not the SAVE_STATE_WARNING
185     if len(caught_warnings) > 1:
186         for w in caught_warnings:
187             if w.category != UserWarning or w.message != SAVE_STATE_WARNING:
188                 warnings.warn(w.message, w.category)
189
190
191 @contextmanager
192 def torch_distributed_zero_first(local_rank: int):
193     """
194     Decorator to make all processes in distributed training wait for each local_master to do something.
195
196     Args:
197         local_rank (:obj:`int`): The rank of the local process.
198
199     """
200     if local_rank not in [-1, 0]:
201         dist.barrier()
202     yield
203     if local_rank == 0:
204         dist.barrier()
205
206
207 class DistributedSamplerWithLoop(DistributedSampler):
208     """
209     Like a :obj:`torch.utils.data.distributed.DistributedSampler` but loops at the end back to the beginning of the
210     shuffled samples to make each process have a round multiple of batch_size samples.
211
212     Args:
213         dataset (:obj:`torch.utils.data.Dataset`):
214             Dataset used for sampling.
215         batch_size (:obj:`int`):
216             The batch size used with this sampler
217         kwargs:
218             All other keyword arguments passed to :obj:`DistributedSampler`.
219
220     """
221
222     def __init__(self, dataset, batch_size, **kwargs):
223         super().__init__(dataset, **kwargs)
224         self.batch_size = batch_size
225
226     def __iter__(self):
227         indices = list(super().__iter__())
228         remainder = 0 if len(indices) % self.batch_size == 0 else self.batch_size - len(indices) % self.batch_size
229         # DistributedSampler already added samples from the beginning to make the number of samples a round multiple
230         # of the world size, so we skip those.
231         start_remainder = 1 if self.rank < len(self.dataset) % self.num_replicas else 0
232         indices += indices[start_remainder : start_remainder + remainder]
233         return iter(indices)
234
235
236 class SequentialDistributedSampler(Sampler):
237     """
238     Distributed Sampler that subsamples indices sequentially, making it easier to collate all results at the end.
239
240     Even though we only use this sampler for eval and predict (no training), which means that the model params won't
241     have to be synced (i.e. will not hang for synchronization even if varied number of forward passes), we still add

```

```

240 extra samples to the sampler to make it evenly divisible (like in `DistributedSampler`) to make it easy to `gather`
241 or `reduce` resulting tensors at the end of the loop.
242 """
243
244 def __init__(self, dataset, num_replicas=None, rank=None, batch_size=None):
245     warnings.warn(
246         "SequentialDistributedSampler is deprecated and will be removed in v5 of Transformers.",
247         FutureWarning,
248     )
249     if num_replicas is None:
250         if not dist.is_available():
251             raise RuntimeError("Requires distributed package to be available")
252         num_replicas = dist.get_world_size()
253     if rank is None:
254         if not dist.is_available():
255             raise RuntimeError("Requires distributed package to be available")
256         rank = dist.get_rank()
257     self.dataset = dataset
258     self.num_replicas = num_replicas
259     self.rank = rank
260     num_samples = len(self.dataset)
261     # Add extra samples to make num_samples a multiple of batch_size if passed
262     if batch_size is not None:
263         self.num_samples = int(math.ceil(num_samples / (batch_size * num_replicas))) * batch_size
264     else:
265         self.num_samples = int(math.ceil(num_samples / num_replicas))
266     self.total_size = self.num_samples * self.num_replicas
267     self.batch_size = batch_size
268
269 def __iter__(self):
270     indices = list(range(len(self.dataset)))
271
272     # add extra samples to make it evenly divisible
273     indices += indices[: (self.total_size - len(indices))]
274     assert (
275         len(indices) == self.total_size
276     ), f"Indices length {len(indices)} and total size {self.total_size} mismatched"
277
278     # subsample
279     indices = indices[self.rank * self.num_samples : (self.rank + 1) * self.num_samples]
280     assert (
281         len(indices) == self.num_samples
282     ), f"Indices length {len(indices)} and sample number {self.num_samples} mismatched"
283
284     return iter(indices)
285
286 def __len__(self):
287     return self.num_samples
288
289
290 def get_tpu_sampler(dataset: torch.utils.data.dataset.Dataset, batch_size: int):
291     if xm.xrt_world_size() <= 1:
292         return RandomSampler(dataset)
293     return DistributedSampler(dataset, num_replicas=xm.xrt_world_size(), rank=xm.get_ordinal())
294
295
296 def nested_new_like(arrays, num_samples, padding_index=-100):
297     """Create the same nested structure as `arrays` with a first dimension always at `num_samples`. """
298     if isinstance(arrays, (list, tuple)):
299         return type(arrays)(nested_new_like(x, num_samples) for x in arrays)
300     return np.full_like(arrays, padding_index, shape=(num_samples, *arrays.shape[1:]))
301
302
303 def expand_like(arrays, new_seq_length, padding_index=-100):
304     """Expand the `arrays` so that the second dimension grows to `new_seq_length`. Uses `padding_index` for padding. """
305     result = np.full_like(arrays, padding_index, shape=(arrays.shape[0], new_seq_length) + arrays.shape[2:])
306     result[:, : arrays.shape[1]] = arrays
307     return result
308
309
310 def nested_truncate(tensors, limit):
311     """Truncate `tensors` at `limit` (even if it's a nested list/tuple of tensors). """
312     if isinstance(tensors, (list, tuple)):
313         return type(tensors)(nested_truncate(t, limit) for t in tensors)
314     return tensors[:limit]
315
316
317 class DistributedTensorGatherer:
318     """
319     A class responsible for properly gathering tensors (or nested list/tuple of tensors) on the CPU by chunks.
320
321     If our dataset has 16 samples with a batch size of 2 on 3 processes and we gather then transfer on CPU at every
322     step, our sampler will generate the following indices:
323
324         :obj:`[0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 0, 1]`
325
326     to get something of size a multiple of 3 (so that each process gets the same dataset length). Then process 0, 1 and
327     2 will be responsible of making predictions for the following samples:
328
329         - P0: :obj:`[0, 1, 2, 3, 4, 5]`
330         - P1: :obj:`[6, 7, 8, 9, 10, 11]`
331         - P2: :obj:`[12, 13, 14, 15, 0, 1]`
332
333     The first batch treated on each process will be
334
335         - P0: :obj:`[0, 1]`
336         - P1: :obj:`[6, 7]`
337         - P2: :obj:`[12, 13]`
338
339     So if we gather at the end of the first batch, we will get a tensor (nested list/tuple of tensor) corresponding to
340     the following indices:
341
342         :obj:`[0, 1, 6, 7, 12, 13]`
343
344     If we directly concatenate our results without taking any precautions, the user will then get the predictions for
345     the indices in this order at the end of the prediction loop:
346
347         :obj:`[0, 1, 6, 7, 12, 13, 2, 3, 8, 9, 14, 15, 4, 5, 10, 11, 0, 1]`
348
349     For some reason, that's not going to roll their boat. This class is there to solve that problem.
350
351     Args:
352
353         world_size (:obj:`int`):
354             The number of processes used in the distributed training.
355         num_samples (:obj:`int`):
356             The number of samples in our dataset.
357         make_multiple_of (:obj:`int`, `optional`):
358             If passed, the class assumes the datasets passed to each process are made to be a multiple of this argument
359             (by adding samples).

```

```

360         padding_index (:obj:'int', 'optional', defaults to -100):
361             The padding index to use if the arrays don't all have the same sequence length.
362     """
363
364     def __init__(self, world_size, num_samples, make_multiple_of=None, padding_index=-100):
365         warnings.warn(
366             "DistributedTensorGatherer is deprecated and will be removed in v5 of Transformers.",
367             FutureWarning,
368         )
369         self.world_size = world_size
370         self.num_samples = num_samples
371         total_size = world_size if make_multiple_of is None else world_size * make_multiple_of
372         self.total_samples = int(np.ceil(num_samples / total_size)) * total_size
373         self.process_length = self.total_samples // world_size
374         self.storage = None
375         self.offsets = None
376         self.padding_index = padding_index
377
378     def add_arrays(self, arrays):
379         """
380         Add :obj:'arrays' to the internal storage, Will initialize the storage to the full size at the first arrays
381         passed so that if we're bound to get an OOM, it happens at the beginning.
382         """
383         if arrays is None:
384             return
385         if self.storage is None:
386             self.storage = nested_new_like(arrays, self.total_samples, padding_index=self.padding_index)
387             self.offsets = list(range(0, self.total_samples, self.process_length))
388
389         slice_len, self.storage = self.nested_set_tensors(self.storage, arrays)
390         for i in range(self.world_size):
391             self.offsets[i] += slice_len
392
393     def nested_set_tensors(self, storage, arrays):
394         if isinstance(arrays, (list, tuple)):
395             result = [self.nested_set_tensors(x, y) for x, y in zip(storage, arrays)]
396             return result[0][0], type(arrays)(r[1] for r in result)
397         assert (
398             arrays.shape[0] % self.world_size == 0
399         ), f"Arrays passed should all have a first dimension multiple of {self.world_size}, found {arrays.shape[0]}."
400
401         slice_len = arrays.shape[0] // self.world_size
402         for i in range(self.world_size):
403             if len(arrays.shape) == 1:
404                 storage[self.offsets[i] : self.offsets[i] + slice_len] = arrays[i * slice_len : (i + 1) * slice_len]
405             else:
406                 # Expand the array on the fly if needed.
407                 if len(storage.shape) > 1 and storage.shape[1] < arrays.shape[1]:
408                     storage = expand_like(storage, arrays.shape[1], padding_index=self.padding_index)
409                 storage[self.offsets[i] : self.offsets[i] + slice_len, : arrays.shape[1]] = arrays[
410                     i * slice_len : (i + 1) * slice_len, :
411                 ]
412         return slice_len, storage
413
414     def finalize(self):
415         """
416         Return the properly gathered arrays and truncate to the number of samples (since the sampler added some extras
417         to get each process a dataset of the same length).
418         """
419         if self.storage is None:
420             return
421         if self.offsets[0] != self.process_length:
422             logger.warning("Not all data has been set. Are you sure you passed all values?")
423         return nested_truncate(self.storage, self.num_samples)
424
425 from typing import Optional
426 from matplotlib import pyplot as plt
427 import numpy as np
428 from numpy.typing import NDArray
429 from pydantic import BaseModel, Extra, Field, validate_arguments
430
431
432 class FCM(BaseModel):
433     """Fuzzy C-means Model"""
434     Attributes:
435         n_clusters (int): The number of clusters to form as well as the number
436         of centroids to generate by the fuzzy C-means.
437         max_iter (int): Maximum number of iterations of the fuzzy C-means
438         algorithm for a single run.
439         m (float): Degree of fuzziness:  $\sum_i \sum_j (u_{ij}^m \|x_i - c_j\|^2) = 0$ .
440         error (float): Relative tolerance with regards to Frobenius norm of
441         the difference
442         in the cluster centers of two consecutive iterations to declare
443         convergence.
444         random_state (Optional[int]): Determines random number generation for
445         centroid initialization.
446         Use an int to make the randomness deterministic.
447         trained (bool): Variable to store whether or not the model has been
448         trained.
449     Returns:
450         FCM: A FCM model.
451     Raises:
452         ReferenceError: If called without the model being trained
453     """
454
455     class Config:
456         extra = Extra.allow
457         arbitrary_types_allowed = True
458
459     n_clusters: int = Field(5, ge=1, le=100)
460     max_iter: int = Field(150, ge=1, le=1000)
461     m: float = Field(2.0, ge=1.0)
462     error: float = Field(1e-5, ge=1e-9)
463     random_state: Optional[int] = None
464     trained: bool = Field(False, const=True)
465
466     @validate_arguments(config=dict(arbitrary_types_allowed=True))
467     def fit(self, X: NDArray) -> None:
468         """Train the fuzzy-c-means model"""
469         Args:
470             X (NDArray): Training instances to cluster
471         """
472         self.rng = np.random.default_rng(self.random_state)
473         n_samples = X.shape[0]
474         self.u = self.rng.uniform(size=(n_samples, self.n_clusters))
475         self.u = self.u / np.tile(
476             self.u.sum(axis=1)[np.newaxis].T, self.n_clusters
477         )
478         for _ in range(self.max_iter):
479             u_old = self.u.copy()

```

```

480         self.centers = FCM_next_centers(X, self.u, self.m)
481         self.u = self.soft_predict(X)
482         # Stopping rule
483         if np.linalg.norm(self.u - u_old) < self.error:
484             break
485         self.trained = True
486
487     @validate_arguments(config=dict(arbitrary_types_allowed=True))
488     def soft_predict(self, X: NDArray) -> NDArray:
489         """Soft predict of FCM
490         Args:
491             X (NDArray): New data to predict.
492         Returns:
493             NDArray: Fuzzy partition array, returned as an array with
494             n_samples rows and n_clusters columns.
495         """
496         temp = FCM_dist(X, self.centers) ** (2 / (self.m - 1))
497         denominator = temp.reshape((X.shape[0], 1, -1)).repeat(
498             temp.shape[-1], axis=1
499         )
500         denominator = temp[:, :, np.newaxis] / denominator
501         return 1 / denominator.sum(2)
502
503     @validate_arguments(config=dict(arbitrary_types_allowed=True))
504     def predict(self, X: NDArray) -> NDArray:
505         """Predict the closest cluster each sample in X belongs to.
506         Args:
507             X (NDArray): New data to predict.
508         Raises:
509             ReferenceError: If it called without the model being trained.
510         Returns:
511             NDArray: Index of the cluster each sample belongs to.
512         """
513         if self.is_trained():
514             X = np.expand_dims(X, axis=0) if len(X.shape) == 1 else X
515             return self.soft_predict(X).argmax(axis=-1)
516         raise ReferenceError(
517             "You need to train the model. Run `.fit()` method to this."
518         )
519
520     def is_trained(self) -> bool:
521         if self.trained:
522             return True
523         return False
524
525     @staticmethod
526     def dist(A: NDArray, B: NDArray) -> NDArray:
527         """Compute the euclidean distance two matrices"""
528         return np.sqrt(np.einsum("ijk->ij", (A[:, None, :] - B) ** 2))
529
530     @staticmethod
531     def next_centers(X: NDArray, u: NDArray, m: float):
532         """Update cluster centers"""
533         um = u ** m
534         return (X.T @ um / np.sum(um, axis=0)).T
535
536     @property
537     def centers(self) -> NDArray:
538         if self.is_trained():
539             return self.centers
540         raise ReferenceError(
541             "You need to train the model. Run `.fit()` method to this."
542         )
543
544     @property
545     def partition_coefficient(self) -> float:
546         """partition coefficient
547         Equation 12a of
548         [this paper](https://doi.org/10.1016/0098-3004(84)90020-7).
549         """
550         if self.is_trained():
551             return np.mean(self.u ** 2)
552         raise ReferenceError(
553             "You need to train the model. Run `.fit()` method to this."
554         )
555
556     @property
557     def partition_coefficient_dev(self) -> float:
558         if self.is_trained():
559             return self.u ** 2
560         raise ReferenceError(
561             "You need to train the model. Run `.fit()` method to this."
562         )
563
564     @property
565     def partition_entropy_coefficient(self):
566         if self.is_trained():
567             return -np.mean(self.u * np.log2(self.u))
568         raise ReferenceError(
569             "You need to train the model. Run `.fit()` method to this."
570         )
571
572     fcm.fit(BigBirdModel.Buffer())
573
574 @dataclass
575 class LabelSmoother:
576     """
577     Adds label-smoothing on a pre-computed output from a Transformers model.
578
579     Args:
580         epsilon (:obj:`float`, 'optional', defaults to 0.1):
581             The label smoothing factor.
582         ignore_index (:obj:`int`, 'optional', defaults to -100):
583             The index in the labels to ignore when computing the loss.
584     """
585     epsilon: float = 0.1
586     ignore_index: int = -100
587
588     def __call__(self, model_output, labels):
589         logits = model_output["logits"] if isinstance(model_output, dict) else model_output[0]
590         log_probs = -torch.nn.functional.log_softmax(logits, dim=-1)
591         if labels.dim() == log_probs.dim() - 1:
592             labels = labels.unsqueeze(-1)
593
594         padding_mask = labels.eq(self.ignore_index)
595         # In case the ignore_index is -100, the gather will fail, so we replace labels by 0. The padding_mask
596         # will ignore them in any case.
597         labels.clamp_min_(0)
598         nll_loss = log_probs.gather(dim=-1, index=labels)
599         # works for fp16 input tensor too, by internally upcasting it to fp32

```

```

600 smoothed_loss = log_probs.sum(dim=-1, keepdim=True, dtype=torch.float32)
601
602 nll_loss.masked_fill_(padding_mask, 0.0)
603 smoothed_loss.masked_fill_(padding_mask, 0.0)
604
605 # Take the mean over the label dimensions, then divide by the number of active elements (i.e. not-padded):
606 num_active_elements = padding_mask.numel() - padding_mask.long().sum()
607 nll_loss = nll_loss.sum() / num_active_elements
608 smoothed_loss = smoothed_loss.sum() / (num_active_elements * log_probs.shape[-1])
609 return (1 - self.epsilon) * nll_loss + self.epsilon * smoothed_loss + torch.var((torch.Tensor(min(1-partition_coefficient_dev[1],partition_coefficient_dev[0])), torch.Tensor(
610
611
612 def get_length_grouped_indices(lengths, batch_size, mega_batch_mult=None, generator=None):
613     """
614     Return a list of indices so that each slice of :obj:`batch_size` consecutive indices correspond to elements of
615     similar lengths. To do this, the indices are:
616
617     - randomly permuted
618     - grouped in mega-batches of size :obj:`mega_batch_mult * batch_size`
619     - sorted by length in each mega-batch
620
621     The result is the concatenation of all mega-batches, with the batch of :obj:`batch_size` containing the element of
622     maximum length placed first, so that an OOM happens sooner rather than later.
623     """
624     # Default for mega_batch_mult: 50 or the number to get 4 megabatches, whichever is smaller.
625     if mega_batch_mult is None:
626         mega_batch_mult = min(len(lengths) // (batch_size * 4), 50)
627     # Just in case, for tiny datasets
628     if mega_batch_mult == 0:
629         mega_batch_mult = 1
630
631     # We need to use torch for the random part as a distributed sampler will set the random seed for torch.
632     indices = torch.randperm(len(lengths), generator=generator)
633     megabatch_size = mega_batch_mult * batch_size
634     megabatches = [indices[i : i + megabatch_size].tolist() for i in range(0, len(lengths), megabatch_size)]
635     megabatches = [list(sorted(megabatch, key=lambda i: lengths[i], reverse=True)) for megabatch in megabatches]
636
637     # The rest is to get the biggest batch first.
638     # Since each megabatch is sorted by descending length, the longest element is the first
639     megabatch_maximums = [lengths[megabatch[0]] for megabatch in megabatches]
640     max_idx = torch.argmax(torch.tensor(megabatch_maximums)).item()
641     # Switch to put the longest element in first position
642     megabatches[0][0], megabatches[max_idx][0] = megabatches[max_idx][0], megabatches[0][0]
643
644     return sum(megabatches, [])
645
646
647 class LengthGroupedSampler(Sampler):
648     """
649     Sampler that samples indices in a way that groups together features of the dataset of roughly the same length while
650     keeping a bit of randomness.
651     """
652
653     def __init__(
654         self,
655         dataset: Dataset,
656         batch_size: int,
657         lengths: Optional[List[int]] = None,
658         model_input_name: Optional[str] = None,
659     ):
660         self.dataset = dataset
661         self.batch_size = batch_size
662         self.model_input_name = model_input_name if model_input_name is not None else "input_ids"
663         if lengths is None:
664             if not isinstance(dataset[0], dict) or self.model_input_name not in dataset[0]:
665                 raise ValueError(
666                     "Can only automatically infer lengths for datasets whose items are dictionaries with an "
667                     f"'{self.model_input_name}' key."
668                 )
669             lengths = [len(feature[self.model_input_name]) for feature in dataset]
670         self.lengths = lengths
671
672     def __len__(self):
673         return len(self.lengths)
674
675     def __iter__(self):
676         indices = get_length_grouped_indices(self.lengths, self.batch_size)
677         return iter(indices)
678
679
680 class DistributedLengthGroupedSampler(DistributedSampler):
681     """
682     Distributed Sampler that samples indices in a way that groups together features of the dataset of roughly the same
683     length while keeping a bit of randomness.
684     """
685     # Copied and adapted from PyTorch DistributedSampler.
686     def __init__(
687         self,
688         dataset: Dataset,
689         batch_size: int,
690         num_replicas: Optional[int] = None,
691         rank: Optional[int] = None,
692         seed: int = 0,
693         drop_last: bool = False,
694         lengths: Optional[List[int]] = None,
695         model_input_name: Optional[str] = None,
696     ):
697         if num_replicas is None:
698             if not dist.is_available():
699                 raise RuntimeError("Requires distributed package to be available")
700             num_replicas = dist.get_world_size()
701         if rank is None:
702             if not dist.is_available():
703                 raise RuntimeError("Requires distributed package to be available")
704             rank = dist.get_rank()
705         self.dataset = dataset
706         self.batch_size = batch_size
707         self.num_replicas = num_replicas
708         self.rank = rank
709         self.epoch = 0
710         self.drop_last = drop_last
711         # If the dataset length is evenly divisible by # of replicas, then there
712         # is no need to drop any data, since the dataset will be split equally.
713         if self.drop_last and len(self.dataset) % self.num_replicas != 0:
714             # Split to nearest available length that is evenly divisible.
715             # This is to ensure each rank receives the same amount of data when
716             # using this Sampler.
717             self.num_samples = math.ceil((len(self.dataset) - self.num_replicas) / self.num_replicas)
718         else:

```

```

720         self.num_samples = math.ceil(len(self.dataset) / self.num_replicas)
721         self.total_size = self.num_samples * self.num_replicas
722         self.seed = seed
723         self.model_input_name = model_input_name if model_input_name is not None else "input_ids"
724
725         if lengths is None:
726             if not isinstance(dataset[0], dict) or self.model_input_name not in dataset[0]:
727                 raise ValueError(
728                     "Can only automatically infer lengths for datasets whose items are dictionaries with an "
729                     f"'{self.model_input_name}' key."
730                 )
731             lengths = [len(feature[self.model_input_name]) for feature in dataset]
732             self.lengths = lengths
733
734     def __iter__(self) -> Iterator:
735         # Deterministically shuffle based on epoch and seed
736         g = torch.Generator()
737         g.manual_seed(self.seed + self.epoch)
738         indices = get_length_grouped_indices(self.lengths, self.batch_size, generator=g)
739
740         if not self.drop_last:
741             # add extra samples to make it evenly divisible
742             indices += indices[: (self.total_size - len(indices))]
743         else:
744             # remove tail of data to make it evenly divisible.
745             indices = indices[: self.total_size]
746         assert len(indices) == self.total_size
747
748         # subsample
749         indices = indices[self.rank : self.total_size : self.num_replicas]
750         assert len(indices) == self.num_samples
751
752         return iter(indices)
753
754 class ShardSampler(Sampler):
755     """
756     Sampler that shards batches between several processes. Dispatches indices batch by batch: on 2 processes with batch
757     size 4, the first two batches are :obj:`[0, 1, 2, 3, 4, 5, 6, 7]` and :obj:`[8, 9, 10, 11, 12, 13, 14, 15]`, which
758     shard into :obj:`[0, 1, 2, 3]` and :obj:`[8, 9, 10, 11]` for GPU-0 and :obj:`[4, 5, 6, 7]` and :obj:`[12, 13, 14,
759     15]` for GPU-1.
760
761     The sampler thus yields :obj:`[0, 1, 2, 3, 8, 9, 10, 11]` on GPU-0 and :obj:`[4, 5, 6, 7, 12, 13, 14, 15]` on
762     GPU-1.
763     """
764
765     def __init__(
766         self,
767         dataset: Dataset,
768         batch_size: int = 1,
769         drop_last: bool = False,
770         num_processes: int = 1,
771         process_index: int = 0,
772     ):
773         self.dataset = dataset
774         self.batch_size = batch_size
775         self.drop_last = drop_last
776         self.num_processes = num_processes
777         self.process_index = process_index
778
779         self.total_batch_size = total_batch_size = batch_size * num_processes
780
781         num_batches = len(dataset) // total_batch_size if drop_last else math.ceil(len(dataset) / total_batch_size)
782         self.total_num_samples = num_batches * total_batch_size
783
784     def __iter__(self):
785         indices = list(range(len(self.dataset)))
786
787         # Add extra samples to make it evenly divisible. While loop is there in the edge case we have a tiny dataset
788         # and it needs to be done several times.
789         while len(indices) < self.total_num_samples:
790             indices += indices[: (self.total_num_samples - len(indices))]
791
792         result = []
793         for batch_start in range(self.batch_size * self.process_index, self.total_num_samples, self.total_batch_size):
794             result += indices[batch_start : batch_start + self.batch_size]
795
796         return iter(result)
797
798     def __len__(self):
799         # Each shard only sees a fraction of total_num_samples.
800         return self.total_num_samples // self.num_processes
801
802 class IterableDatasetShard(IterableDataset):
803     """
804     Wraps a PyTorch :obj:`IterableDataset` to generate samples for one of the processes only. Instances of this class
805     will always yield a number of samples that is a round multiple of the actual batch size (which is :obj:`batch_size
806     x num_processes`). Depending on the value of the :obj:`drop_last` attribute, it will either stop the iteration at
807     the first batch that would be too small or loop with indices from the beginning.
808
809     On two processes with an iterable dataset yielding of :obj:`[0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11]` with a batch
810     size of 2:
811
812     - the shard on process 0 will yield :obj:`[0, 1, 4, 5, 8, 9]` so will see batches :obj:`[0, 1]`, :obj:`[4, 5]`,
813       :obj:`[8, 9]`
814     - the shard on process 1 will yield :obj:`[2, 3, 6, 7, 10, 11]` so will see batches :obj:`[2, 3]`, :obj:`[6, 7]`,
815       :obj:`[10, 11]`
816
817     .. warning:
818
819         If your IterableDataset implements some randomization that needs to be applied the same way on all processes
820         (for instance, a shuffling), you should use a :obj:`torch.Generator` in a :obj:`generator` attribute of the
821         :obj:`dataset` to generate your random numbers and call the
822         :meth:`~transformers.trainer_pt_utils.IterableDatasetShard.set_epoch` method of this object. It will set the
823         seed of this :obj:`generator` to :obj:`seed + epoch` on all processes before starting the iteration.
824         Alternatively, you can also implement a :obj:`set_epoch()` method in your iterable dataset to deal with this.
825
826     Args:
827         dataset (:obj:`torch.utils.data.dataset.IterableDataset`):
828             The batch sampler to split in several shards.
829         batch_size (:obj:`int`, 'optional', defaults to 1):
830             The size of the batches per shard.
831         drop_last (:obj:`bool`, 'optional', defaults to :obj:`False`):
832             Whether or not to drop the last incomplete batch or complete the last batches by using the samples from the
833             beginning.
834         num_processes (:obj:`int`, 'optional', defaults to 1):
835             The number of processes running concurrently.
836         process_index (:obj:`int`, 'optional', defaults to 0):

```

```

840         The index of the current process.
841     seed (:obj:`int`, 'optional', defaults to 0):
842         A random seed that will be used for the random number generation in
843         :meth:`~transformers.trainer_pt_utils.IterableDatasetShard.set_epoch`.
844     """
845
846     def __init__(
847         self,
848         dataset: IterableDataset,
849         batch_size: int = 1,
850         drop_last: bool = False,
851         num_processes: int = 1,
852         process_index: int = 0,
853         seed: int = 0,
854     ):
855         self.dataset = dataset
856         self.batch_size = batch_size
857         self.drop_last = drop_last
858         self.num_processes = num_processes
859         self.process_index = process_index
860         self.seed = seed
861         self.epoch = 0
862         self.num_examples = 0
863
864     def set_epoch(self, epoch):
865         self.epoch = epoch
866         if hasattr(self.dataset, "set_epoch"):
867             self.dataset.set_epoch(epoch)
868
869     def __iter__(self):
870         self.num_examples = 0
871         if (
872             not hasattr(self.dataset, "set_epoch")
873             and hasattr(self.dataset, "generator")
874             and isinstance(self.dataset.generator, torch.Generator)
875         ):
876             self.dataset.generator.manual_seed(self.seed + self.epoch)
877         real_batch_size = self.batch_size * self.num_processes
878         process_slice = range(self.process_index * self.batch_size, (self.process_index + 1) * self.batch_size)
879
880         first_batch = None
881         current_batch = []
882         for element in self.dataset:
883             self.num_examples += 1
884             current_batch.append(element)
885             # Wait to have a full batch before yielding elements.
886             if len(current_batch) == real_batch_size:
887                 for i in process_slice:
888                     yield current_batch[i]
889                 if first_batch is None:
890                     first_batch = current_batch.copy()
891                 current_batch = []
892
893             # Finished if drop_last is True, otherwise complete the last batch with elements from the beginning.
894             if not self.drop_last and len(current_batch) > 0:
895                 if first_batch is None:
896                     first_batch = current_batch.copy()
897                 while len(current_batch) < real_batch_size:
898                     current_batch += first_batch
899                 for i in process_slice:
900                     yield current_batch[i]
901
902 # In order to keep `trainer.py` compact and easy to understand, place any secondary PT Trainer
903 # helper methods here
904
905
906 def _get_learning_rate(self):
907     if self.deepspeed:
908         # with deepspeed's fp16 and dynamic loss scale enabled the optimizer/scheduler steps may
909         # not run for the first few dozen steps while loss scale is too large, and thus during
910         # that time `get_last_lr` will fail if called during that warm up stage, so work around it:
911         try:
912             last_lr = self.lr_scheduler.get_last_lr()[0]
913         except AssertionError as e:
914             if "need to call step" in str(e):
915                 logger.warning("Tried to get lr value before scheduler/optimizer started stepping, returning lr=0")
916                 last_lr = 0
917             else:
918                 raise
919     else:
920         last_lr = (
921             # backward compatibility for pytorch schedulers
922             self.lr_scheduler.get_last_lr()[0]
923             if version.parse(torch.__version__) >= version.parse("1.4")
924             else self.lr_scheduler.get_lr()[0]
925         )
926     return last_lr
927
928
929 def _secs2timedelta(secs):
930     """
931     convert seconds to hh:mm:ss.msec, msec rounded to 2 decimals
932     """
933     msec = int(abs(secs - int(secs)) * 100)
934     return f"[datetime.timedelta(seconds=int(secs))].(msec:02d)"
935
936
937 def metrics_format(self, metrics: Dict[str, float]) -> Dict[str, float]:
938     """
939     Reformat Trainer metrics values to a human-readable format
940
941     Args:
942         metrics (:obj:`Dict[str, float]`):
943             The metrics returned from train/evaluate/predict
944
945     Returns:
946         metrics (:obj:`Dict[str, float]`): The reformatted metrics
947     """
948     metrics_copy = metrics.copy()
949     for k, v in metrics_copy.items():
950         if "mem" in k:
951             metrics_copy[k] = f"{v >> 20 }MB"
952         elif "runtime" in k:
953             metrics_copy[k] = _secs2timedelta(v)
954         elif k == "total_flos":
955             metrics_copy[k] = f"{int(v) >> 30 }GF"
956         elif type(metrics_copy[k]) == float:

```



```

960         metrics_copy[k] = round(v, 4)
961
962     return metrics_copy
963
964
965 def log_metrics(self, split, metrics):
966     """
967     Log metrics in a specially formatted way
968
969     Under distributed environment this is done only for a process with rank 0.
970
971     Args:
972         split (:obj:`str`):
973             Mode/split name: one of ``train``, ``eval``, ``test``
974         metrics (:obj:`Dict[str, float]`):
975             The metrics returned from train/evaluate/predictmetrics: metrics dict
976
977     Notes on memory reports:
978
979     In order to get memory usage report you need to install ``psutil``. You can do that with ``pip install psutil``.
980
981     Now when this method is run, you will see a report that will include ::
982
983         init_mem_cpu_alloc_delta = 1301MB
984         init_mem_cpu_peaked_delta = 154MB
985         init_mem_gpu_alloc_delta = 230MB
986         init_mem_gpu_peaked_delta = 0MB
987         train_mem_cpu_alloc_delta = 1345MB
988         train_mem_cpu_peaked_delta = 0MB
989         train_mem_gpu_alloc_delta = 693MB
990         train_mem_gpu_peaked_delta = 7MB
991
992     **Understanding the reports:**
993
994     - the first segment, e.g., ``train``, tells you which stage the metrics are for. Reports starting with ``init``
995       will be added to the first stage that gets run. So that if only evaluation is run, the memory usage for the
996       ``init`` will be reported along with the ``eval`` metrics.
997     - the third segment, is either ``cpu`` or ``gpu``, tells you whether it's the general RAM or the gpu0 memory
998       metric.
999     - ``*_alloc_delta`` - is the difference in the used/allocated memory counter between the end and the start of the
1000       stage - it can be negative if a function released more memory than it allocated.
1001     - ``*_peaked_delta`` - is any extra memory that was consumed and then freed - relative to the current allocated
1002       memory counter - it is never negative. When you look at the metrics of any stage you add up ``*_alloc_delta`` +
1003       ``*_peaked_delta`` and you know how much memory was needed to complete that stage.
1004
1005     The reporting happens only for process of rank 0 and gpu 0 (if there is a gpu). Typically this is enough since the
1006     main process does the bulk of work, but it could be not quite so if model parallel is used and then other GPUs may
1007     use a different amount of gpu memory. This is also not the same under DataParallel where gpu0 may require much more
1008     memory than the rest since it stores the gradient and optimizer states for all participating GPUS. Perhaps in the
1009     future these reports will evolve to measure those too.
1010
1011     The CPU RAM metric measures RSS (Resident Set Size) includes both the memory which is unique to the process and the
1012     memory shared with other processes. It is important to note that it does not include swapped out memory, so the
1013     reports could be imprecise.
1014
1015     The CPU peak memory is measured using a sampling thread. Due to python's GIL it may miss some of the peak memory if
1016     that thread didn't get a chance to run when the highest memory was used. Therefore this report can be less than
1017     reality. Using ``tracemalloc`` would have reported the exact peak memory, but it doesn't report memory allocations
1018     outside of python. So if some C++ CUDA extension allocated its own memory it won't be reported. And therefore it
1019     was dropped in favor of the memory sampling approach, which reads the current process memory usage.
1020
1021     The GPU allocated and peak memory reporting is done with ``torch.cuda.memory_allocated()`` and
1022     ``torch.cuda.max_memory_allocated()``. This metric reports only "deltas" for pytorch-specific allocations, as
1023     ``torch.cuda`` memory management system doesn't track any memory allocated outside of pytorch. For example, the
1024     very first cuda call typically loads CUDA kernels, which may take from 0.5 to 2GB of GPU memory.
1025
1026     Note that this tracker doesn't account for memory allocations outside of :class:`~transformers.Trainer`'s
1027     ``__init``, ``train``, ``evaluate`` and ``predict`` calls.
1028
1029     Because ``evaluation`` calls may happen during ``train``, we can't handle nested invocations because
1030     ``torch.cuda.max_memory_allocated`` is a single counter, so if it gets reset by a nested eval call, ``train``'s
1031     tracker will report incorrect info. If this 'pytorch issue <https://github.com/pytorch/pytorch/issues/16266>'
1032     gets resolved it will be possible to change this class to be re-entrant. Until then we will only track the outer
1033     level of ``train``, ``evaluate`` and ``predict`` methods. Which means that if ``eval`` is called during ``train``,
1034     it's the latter that will account for its memory usage and that of the former.
1035
1036     This also means that if any other tool that is used along the :class:`~transformers.Trainer` calls
1037     ``torch.cuda.reset_peak_memory_stats``, the gpu peak memory stats could be invalid. And the
1038     :class:`~transformers.Trainer` will disrupt the normal behavior of any such tools that rely on calling
1039     ``torch.cuda.reset_peak_memory_stats`` themselves.
1040
1041     For best performance you may want to consider turning the memory profiling off for production runs.
1042
1043     """
1044     if not self.is_world_process_zero():
1045         return
1046
1047     logger.info(f"***** {split} metrics *****")
1048     metrics_formatted = self.metrics_format(metrics)
1049     k_width = max(len(str(x)) for x in metrics_formatted.keys())
1050     v_width = max(len(str(x)) for x in metrics_formatted.values())
1051     for key in sorted(metrics_formatted.keys()):
1052         logger.info(f"({key: <{k_width}}) = {metrics_formatted[key]:>{v_width}}")
1053
1054 def save_metrics(self, split, metrics, combined=True):
1055     """
1056     Save metrics into a json file for that split, e.g. ``train_results.json``.
1057
1058     Under distributed environment this is done only for a process with rank 0.
1059
1060     Args:
1061         split (:obj:`str`):
1062             Mode/split name: one of ``train``, ``eval``, ``test``, ``all``
1063         metrics (:obj:`Dict[str, float]`):
1064             The metrics returned from train/evaluate/predict
1065         combined (:obj:`bool`, `optional`, defaults to :obj:`True`):
1066             Creates combined metrics by updating ``all_results.json`` with metrics of this call
1067
1068     To understand the metrics please read the docstring of :meth:`~transformers.Trainer.log_metrics`. The only
1069     difference is that raw unformatted numbers are saved in the current method.
1070
1071     """
1072     if not self.is_world_process_zero():
1073         return
1074
1075     path = os.path.join(self.args.output_dir, f"{split}_results.json")
1076     with open(path, "w") as f:
1077         json.dump(metrics, f, indent=4, sort_keys=True)
1078
1079     if combined:

```

```

1080     path = os.path.join(self.args.output_dir, "all_results.json")
1081     if os.path.exists(path):
1082         with open(path, "r") as f:
1083             all_metrics = json.load(f)
1084     else:
1085         all_metrics = {}
1086
1087     all_metrics.update(metrics)
1088     with open(path, "w") as f:
1089         json.dump(all_metrics, f, indent=4, sort_keys=True)
1090
1091 def save_state(self):
1092     """
1093     Saves the Trainer state, since Trainer.save_model saves only the tokenizer with the model
1094
1095     Under distributed environment this is done only for a process with rank 0.
1096     """
1097     if not self.is_world_process_zero():
1098         return
1099
1100     path = os.path.join(self.args.output_dir, "trainer_state.json")
1101     self.state.save_to_json(path)
1102
1103 def get_parameter_names(model, forbidden_layer_types):
1104     """
1105     Returns the names of the model parameters that are not inside a forbidden layer.
1106     """
1107     result = []
1108     for name, child in model.named_children():
1109         result += [
1110             f"{name}.{n}"
1111             for n in get_parameter_names(child, forbidden_layer_types)
1112             if not isinstance(child, tuple(forbidden_layer_types))
1113         ]
1114     # Add model specific parameters (defined with nn.Parameter) since they are not in any child.
1115     result += list(model._parameters.keys())
1116     return result
1117
1118 if is_sagemaker_mp_enabled():
1119     import smdistributed.modelparallel.torch as smp
1120
1121     @smp.step()
1122     def smp_forward_backward(model, inputs, gradient_accumulation_steps=1, scaler=None):
1123         with torch.cuda.amp.autocast(enabled=(scaler is not None)):
1124             outputs = model(**inputs)
1125
1126             loss = outputs["loss"] if isinstance(outputs, dict) else outputs[0]
1127             loss /= gradient_accumulation_steps
1128             if scaler is not None:
1129                 loss = scaler.scale(loss).squeeze()
1130
1131             model.backward(loss)
1132             return loss
1133
1134     @smp.step()
1135     def smp_forward_only(model, inputs):
1136         return model(**inputs)
1137
1138     def smp_gather(tensor):
1139         if isinstance(tensor, (list, tuple)):
1140             return type(tensor)(smp_gather(t) for t in tensor)
1141         elif isinstance(tensor, dict):
1142             return type(tensor)({k: smp_gather(v) for k, v in tensor.items()})
1143         elif not isinstance(tensor, torch.Tensor):
1144             raise TypeError(
1145                 f"Can't gather the values of type {type(tensor)}, only of nested list/tuple/dicts of tensors."
1146             )
1147         all_tensors = smp.allgather(tensor, smp.CommGroup.DP_GROUP)
1148         return torch.cat([t.cpu() for t in all_tensors], dim=0)
1149
1150     def smp_nested_concat(tensor):
1151         if isinstance(tensor, (list, tuple)):
1152             return type(tensor)(smp_nested_concat(t) for t in tensor)
1153         elif isinstance(tensor, dict):
1154             return type(tensor)({k: smp_nested_concat(v) for k, v in tensor.items()})
1155         # It doesn't seem possible to check here if `tensor` is a StepOutput because StepOutput lives in `smp.step`
1156         # which is also the name of the decorator so Python is confused.
1157         return tensor.concat().detach().cpu()

```