

```
1: """ Core mathematics for Additive Quantization (AQ): initialization, reconstruction and beam search """
2: import random
3: from typing import List, Optional, Tuple, Union
4:
5: import torch
6: import torch.nn as nn
7: import torch.nn.functional as F
8: from torch.utils.checkpoint import checkpoint
9: from tqdm.auto import trange
10:
11: from src.kmeans import find_nearest_cluster, fit_faiss_kmeans, fit_kmeans, fit_kmeans_ld
12: from src.utils import ellipsis, maybe_script
13:
14:
15: class QuantizedLinear(nn.Module):
16:     def __init__(self, quantized_weight, bias: Optional[nn.Parameter]):
17:         super().__init__()
18:         self.out_features, self.in_features = quantized_weight.out_features, quantized_weight.in_features
19:         self.quantized_weight = quantized_weight
20:         self.bias = bias
21:         self.use_checkpoint = False
22:
23:     def _forward(self, input: torch.Tensor):
24:         return F.linear(input, self.quantized_weight(), self.bias)
25:
26:     def forward(self, input: torch.Tensor):
27:         if self.use_checkpoint and torch.is_grad_enabled():
28:             return checkpoint(
29:                 self._forward, input, use_reentrant=False, preserve_rng_state=False, deterministic_check="none"
30:             )
31:         return self._forward(input)
32:
33:
34: class QuantizedWeight(nn.Module):
35:     EPS = 1e-9
36:
37:     def __init__(
38:         self,
39:         *,
40:         XTX: torch.Tensor,
41:         reference_weight: torch.Tensor,
42:         in_group_size: int,
43:         out_group_size: int,
44:         num_codebooks: int,
45:         nbits_per_codebook: int = 8,
46:         codebook_value_nbits: int = 16,
47:         codebook_value_num_groups: int = 1,
48:         scale_nbits: int = 0,
49:         scale_in_group_size: Optional[int] = None,
50:         scale_out_group_size: Optional[int] = None,
51:         straight_through_gradient: Optional[bool] = None,
52:         **init_kwargs,
53:     ):
54:         super().__init__()
55:         self.out_features, self.in_features = reference_weight.shape
56:         assert self.in_features % in_group_size == 0
57:         assert self.out_features % out_group_size == 0
58:
59:         self.out_group_size, self.in_group_size = out_group_size, in_group_size
60:         self.num_codebooks = num_codebooks
61:         self.nbits_per_codebook = nbits_per_codebook
62:         self.codebook_size = codebook_size = 2**nbits_per_codebook
63:         self.codebook_value_nbits = codebook_value_nbits
64:         self.codebook_value_num_groups = codebook_value_num_groups
65:         self.codebook_value_clusters = None
66:
67:         if scale_in_group_size is None:
68:             scale_in_group_size = in_group_size
69:         if scale_out_group_size is None:
70:             scale_out_group_size = out_group_size
71:         assert scale_out_group_size % out_group_size == 0
72:         assert scale_in_group_size % in_group_size == 0
```

```

73:         self.scale_in_group_size_factor = scale_in_group_size // in_group_size
74:         self.scale_out_group_size_factor = scale_out_group_size // out_group_size
75:
76:
77:         self.scales = self.scales_clusters = self.scales_indices = None
78:         if straight_through_gradient is None and scale_nbits > 0:
79:             straight_through_gradient = scale_nbits >= 6
80:         self.straight_through_gradient = straight_through_gradient
81:         self.scale_nbits = scale_nbits
82:
83:         with torch.no_grad():
84:             weight_groupwise_for_scales = reference_weight.reshape(
85:                 self.out_features // scale_out_group_size, scale_out_group_size,
86:                 self.in_features // scale_in_group_size, scale_in_group_size
87:             ).swapaxes(1, 2) # [num_out_groups, num_in_groups, out_group_size, in_group_size]
88:
89:             if scale_nbits > 0:
90:                 scales = weight_groupwise_for_scales.norm(dim=(2, 3), keepdim=True) + self.EPS
91:             else:
92:                 scales = weight_groupwise_for_scales.flatten(1, -1).norm(dim=-1).view(-1, 1, 1, 1) + self.EPS
93:             # shape [num_out_groups, num_in_groups, 1, 1] if scale_nbits > 0 else [num_out_groups, num_in_groups, 1, 1]
94:             del weight_groupwise_for_scales
95:
96:
97:             weight_groupwise = reference_weight.reshape(
98:                 self.out_features // out_group_size, out_group_size, self.in_features // in_group_size, in_group_size
99:             ).swapaxes(1, 2) # [num_out_groups, num_in_groups, out_group_size, in_group_size]
100:
101:
102:             self.scales_are_lossless = scale_nbits == 0 or scale_nbits >= 16 or (2**scale_nbits > scales.shape[1])
103:             if self.scales_are_lossless or self.straight_through_gradient:
104:                 # ^-- this checks if scales can be preserved losslessly
105:                 self.scales = nn.Parameter(scales, requires_grad=True)
106:             else:
107:                 scales_clusters, scales_indices, _ = fit_kmeans_1d(scales.flatten(1, -1), k=2**scale_nbits)
108:                 self.scales_clusters = nn.Parameter(scales_clusters, requires_grad=True)
109:                 self.scales_indices = nn.Parameter(scales_indices, requires_grad=False)
110:
111:             weight_for_init = (weight_groupwise / self.get_scales()).swapaxes(1, 2).reshape_as(reference_weight)
112:             del weight_groupwise
113:
114:             codes, codebooks = init_aq_kmeans(
115:                 weight_for_init,
116:                 num_codebooks=num_codebooks,
117:                 out_group_size=out_group_size,
118:                 in_group_size=in_group_size,
119:                 codebook_size=self.codebook_size,
120:                 **init_kwargs,
121:             )
122:
123:             self.codebooks = nn.Parameter(
124:                 codebooks, requires_grad=True
125:             ) # [num_codebooks, codebook_size, out_group_size, in_group_size]
126:             self.codes = nn.Parameter(codes, requires_grad=False) # [num_out_groups, num_in_groups, num_codebooks]
127:
128:             def get_codebooks(self) -> torch.Tensor:
129:                 """Get quantization codebooks or reconstruct them from second level quantization (see codebook_values_nbits)"""
130:                 if self.codebook_value_nbits >= 16:
131:                     return self.codebooks
132:                 elif 0 < self.codebook_value_nbits < 16:
133:                     with torch.no_grad():
134:                         codebooks_dimshuffle = (
135:                             self.codebooks.reshape(
136:                                 self.num_codebooks,
137:                                 self.codebook_value_num_groups,
138:                                 self.codebook_size // self.codebook_value_num_groups,
139:                                 self.out_group_size,

```

```

aq.py Fri Mar 22 16:01:31 2024 3
140:         self.in_group_size,
141:     )
142:     .permute(0, 1, 3, 4, 2)
143:     .flatten(0, -2)
144: )
145: self.codebook_value_clusters, _unused, reconstructed_codebooks_dimshuffle = fit_kmeans_1d(
146:     codebooks_dimshuffle,
147:     k=2**self.codebook_value_nbits,
148:     initial_clusters=self.codebook_value_clusters,
149: )
150: reconstructed_codebooks = (
151:     reconstructed_codebooks_dimshuffle.view(
152:         self.num_codebooks,
153:         self.codebook_value_num_groups,
154:         self.out_group_size,
155:         self.in_group_size,
156:         self.codebook_size // self.codebook_value_num_groups,
157:     )
158:     .permute(0, 1, 4, 2, 3)
159:     .reshape_as(self.codebooks)
160: )
161: if torch.is_grad_enabled():
162:     reconstructed_codebooks = reconstructed_codebooks + (self.codebooks - self.codebooks.detach())
163:     return reconstructed_codebooks
164:     raise NotImplementedError(f"{self.codebook_value_nbits}-bit codebook values are not supported")
165:
166: def get_scales(self) -> torch.Tensor:
167:     """Get per-channel or per-group quantization scales or reconstruct those scales based on scales_nbits"""
168:     if self.scale_nbits == 0 or self.scales_are_lossless:
169:         scales = self.scales # scales are not quantized or the quantization is lossless
170:     elif self.straight_through_gradient:
171:         with torch.no_grad():
172:             self.scales_clusters, _, dequantized_scales = fit_kmeans_1d(
173:                 self.scales.flatten(1, -1), k=2**self.scale_nbits, initial_clusters=self.scales_clusters
174:             )
175:             dequantized_scales = dequantized_scales.reshape_as(self.scales)
176:             if torch.is_grad_enabled() and self.scales.requires_grad:
177:                 dequantized_scales = dequantized_scales + (self.scales - self.scales.detach())
178:             scales = dequantized_scales
179:     else: # train scale codebook only
180:         scales = self.scales_clusters.gather(1, self.scales_indices)[: , :, None, None]
181:
182:         if scales.numel() != scales.shape[0]: # group-wise scales (i.e. not 1d scales)
183:             assert scales.ndim == 4
184:             assert scales.shape[2] == scales.shape[3] == 1
185:             # if replicate each scale several times
186:             scales = scales[: , None, :, None, :, :].tile(
187:                 1, self.scale_out_group_size_factor, 1, self.scale_in_group_size_factor, 1, 1
188:             ).flatten(2, 3).flatten(0, 1)
189:         return scales
190:
191: def forward(self, selection: Union[slice, ellipsis, torch.Tensor] = ...):
192:     """
193:     Differentably reconstruct the weight (or parts thereof) from compressed components
194:     :param selection: By default, reconstruct the entire weight. If selection is specified, this method will instead
195:     reconstruct a portion of weight for the corresponding output dimensions (used for parallelism).
196:     The indices / slices must correspond to output channels (if out_group_size==1) or groups (if > 1).
197:     Formally, the indices must be in range [ 0 , self.out_features // self.out_group_size )
198:
199:     """
200:     weight = _dequantize_weight(self.codes[selection], self.get_codebooks(), self.get_scales()[selection])
201:     return weight
202:
203: @torch.no_grad()
204: def beam_search_update_codes_

```

```

205:         self,
206:         XTX: torch.Tensor,
207:         reference_weight: torch.Tensor,
208:         *,
209:         selection: Union[slice, ellipsis, torch.LongTensor] = ...,
210:         **kwargs,
211:     ) -> torch:
212:         """
213:         Update self.codes in-place via beam search so as to minimize squared errors. Return the updated codes.
214:         :param XTX: pairwise products of input features matmul(X.transpose(), X), shape: [in_features, in_features]
215:         :note: if XTX is divided by dataset size, this function will return *mean* squared error
216:         :param reference_weight: original weight matrix that is being quantized, shape: [out_features, in_features]
217:         :note: if selection is specified, reference_weight must instead be [num_selected_out_features, in_features]
218:         :param selection: By default, this function updates all codes, If selection specified, it will instead
219:             update only the codes for a portion of output dimensions (used for parallelism).
220:             The indices / slices must correspond to output channels (if out_group_size==1) or groups (if > 1).
221:             Formally, the indices must be in range [ 0 , self.out_features // self.out_group_size )
222:         :param beam_size: consider up to this many best encoding combinations (this param is passed through via kwargs)
223:         :param kwargs: any additional keyword arguments are forwarded to beam_search_optimal_codes function
224:         :returns: the updated codes
225:         """
226:         if self.scale_out_group_size_factor != 1:
227:             assert selection == ..., "todo implement selection with scale_out_group_size"
228:         self.codes[selection] = beam_search_optimal_codes(
229:             XTX=XTX,
230:             reference_weight=reference_weight,
231:             codebooks=self.get_codebooks(),
232:             prev_codes=self.codes[selection],
233:             scales=self.get_scales()[selection],
234:             **kwargs,
235:         )
236:         return self.codes[selection]
237:
238:     def estimate_nbits_per_parameter(self) -> float:
239:         """Calculate the effective number of bits per original matrix parameters"""
240:         num_parameters = self.out_features * self.in_features
241:         group_size = self.out_group_size * self.in_group_size
242:         num_out_groups = self.out_features // self.out_group_size
243:         num_in_groups = self.in_features // self.in_group_size
244:
245:         matrix_store = num_parameters // group_size * self.num_codebooks * self.nbits_per_codebook
246:
247:         codebooks_store = self.num_codebooks * self.codebook_size * group_size * self.codebook_value_nbits
248:         if self.codebook_value_nbits < 16:
249:             codebooks_store += (
250:                 2**self.codebook_value_nbits * self.num_codebooks * self.codebook_value_num_groups * group_size * 16
251:             )
252:
253:         if self.scale_nbits >= 16 or 2**self.scale_nbits >= num_in_groups: # group-wise scales in 16 bit
254:             scale_num_out_groups = num_out_groups // self.scale_out_group_size_factor
255:             scale_num_in_groups = num_in_groups // self.scale_in_group_size_factor
256:             scale_store = self.scale_nbits * scale_num_out_groups * scale_num_in_groups
257:         elif 0 < self.scale_nbits < 16: # use scale quantization codebooks
258:             scale_store = self.scale_nbits * num_out_groups * num_in_groups
259:             scale_store += num_out_groups * 2**self.scale_nbits * 16
260:         elif self.scale_nbits == 0: # no group-wise scales; use global 1d scales instead
261:             scale_store = num_out_groups * 16
262:         else:
263:             assert False
264:
265:         return (matrix_store + codebooks_store + scale_store) / num_parameters

```

```

aq.py           Fri Mar 22 16:01:31 2024           5

266:
267:     def extra_repr(self) -> str:
268:         return f"{self.out_features=}, {self.in_features=}, bits_per_parameter={self.estimate_
nbits_per_parameter()}"
269:
270:
271: @torch.inference_mode()
272: def beam_search_optimal_codes(
273:     *,
274:     XTX: torch.Tensor,
275:     reference_weight: torch.Tensor,
276:     codebooks: torch.Tensor,
277:     prev_codes: torch.IntTensor,
278:     scales: Optional[torch.Tensor],
279:     beam_size: int,
280:     dim_rng: Optional[random.Random] = None,
281:     code_penalties: Optional[torch.Tensor] = None,
282:     verbose: bool,
283: ):
284:     """
285:     :param XTX: pairwise products of input features matmul(X.transpose(), X), shape: [in_featu
res, in_features]
286:     :note: if XTX is divided by dataset size, this function will return *mean* squared error
287:     :param reference_weight: original weight matrix that is being quantized, shape: [out_featu
res, in_features]
288:     :param codebooks: look-up tables of codes, shape: [num_codebooks, codebook_size, out_group
_size, in_group_size]
289:     :param prev_codes: previous-best integer weight codes, shape: [num_out_groups, num_in_grou
ps, num_codebooks]
290:     :param scales: weight will be multiplied by this factor, shape = [num_out_groups, num_in_g
roups or 1, 1, 1]
291:     :param dim_rng: a source of randomness to (optionally) shuffle the order in which the beam
search runs
292:     None = update dimensions and codebooks in their natural order (0, 1, ..., n)
293:     random.Random(optional_seed) = shuffle dimensions at random, optionally using the specif
ied seed
294:
295:     :param beam_size: consider up to this many best encoding combinations
296:     :param code_penalties: a pytorch float tensor of shape [num_codebooks, codebook_size]
297:     Penalize the beam search objective by code_penalties[m][i] for every use of ith code f
rom mth codebook
298:     :param verbose: if True, draw a progressbar and periodically print best loss
299:     :return: best quantization codes found, same shape as prev_codes
300:
301:     :intuition: the beam search needs to produce weight codes that minimize MSE error
302:     - the codes are of shape [out_features / out_group_size, in_features / in_group_size, num_
codebooks]
303:
304:     Out of those three dimensions, out_features is "independent", i.e. changing code in
305:     one output feature does not increase the MSE error for another feature. Therefore,
306:     beam search for different output features can run in independently in parallel.
307:
308:     Neither (in_features / in_group_size) nor (num_codebooks) dimension are independent:
309:     - changing the encoding for one feature can compensate the error from encoding another, OB
C-style
310:     - for a single weight group, changing code in one codebook can affect the optimal choice i
n another codebook
311:     Therefore, beam search must go in a double loop over (in_features/in_group_size) and (num_
codebooks) dimensions
312:
313:     This leaves one choice: which dimension used for outer loop, and which one goes is in the
inner loop?
314:     Due to the nature of beam search, interactions between dimensions of inner loop will be ex
plored better.
315:     We chose to use (in_features/in_group_size) in the outer loop and (num_codebooks) for the
inner loop.
316:     This is based on an intuition from GPTQ: you can get decent performance by quantizing each
input unit ...
317:     ... greedily --- GPTQ does not change quantizations for previously quantized features and
works fine.
318:     Therefore, we believe that we can also use a greedy approach to compensate error between i
nput features.
319:     In turn, we believe that the codes used to encode the same weights (additively) are more i
nter-dependent.
320:     This should be treated as an educated guess with no proof and no ablation (as of the time

```

```

of writing).
321:
322:     """
323:     num_out_groups, num_in_groups, num_codebooks = prev_codes.shape
324:     num_codebooks, codebook_size, out_group_size, in_group_size = codebooks.shape
325:     in_features = num_in_groups * in_group_size
326:     out_features = num_out_groups * out_group_size
327:     assert reference_weight.shape == (out_features, in_features)
328:     prev_weight = _dequantize_weight(prev_codes, codebooks, scales)
329:
330:     # initialize all beam codes as previous codes - so they can be updated during beam search
331:     beam_codes = prev_codes.unsqueeze(0)
332:     # beam_codes shape: [current beam_size, num_out_groups, num_in_groups, num_codebooks], ini
    tial beam_size = 1
333:     beam_weights = prev_weight.unsqueeze(0)
334:     # beam_weights shape: [current beam_size, out_features, in_features], initial beam size =
1
335:
336:     beam_losses = (
337:         _channelwise_squared_error(CTX, prev_weight, reference_weight)
338:         .reshape(1, num_out_groups, out_group_size)
339:         .sum(-1)
340:     )
341:     # beam_losses shape: [current beam_size, num_out_groups], initial beam_size = 1
342:     if code_penalties is not None:
343:         # Compute counts for each code in each codebook, initialize regularizer
344:         codebook_ids = torch.arange(num_codebooks, device=beam_losses.device).view(1, 1, 1, -1
)
345:         per_channel_regularizers = code_penalties[codebook_ids, beam_codes].sum(dim=(2, 3)) #
[beam_size, num_out_groups]
346:         del codebook_ids
347:         beam_losses += beam_losses + per_channel_regularizers
348:
349:     if verbose:
350:         progressbar = trange(num_in_groups * num_codebooks)
351:
352:     def _make_range(n: int) -> list:
353:         seq = list(range(n))
354:         if dim_rng is not None:
355:             dim_rng.shuffle(seq)
356:         return seq
357:
358:     for input_group_index in _make_range(num_in_groups):
359:         for codebook_index in _make_range(num_codebooks):
360:             ### part 1: compute losses for every possible candidate for one given codebook and
input group.
361:             # Currently, we compute errors for all output features in parallel in a vectorized
fashion.
362:             best_losses, best_indices = _beam_search_squared_errors(
363:                 CTX=CTX,
364:                 reference_weight=reference_weight,
365:                 codebooks=codebooks,
366:                 scales=scales,
367:                 beam_losses=beam_losses,
368:                 beam_codes=beam_codes,
369:                 beam_weights=beam_weights,
370:                 input_group_index=input_group_index,
371:                 codebook_index=codebook_index,
372:                 k_best=beam_size,
373:                 code_penalties=code_penalties,
374:             ) # [current beam_size, codebook_size, num_out_groups]
375:
376:             # part 2: select beam_size new best codes and re-arrange beam to account for the f
act that ...
377:             # ... sometimes two or more top candidates originate from the same source in previ
ous beam
378:             beam_codes, beam_weights, beam_losses = _beam_search_select_best(
379:                 beam_codes=beam_codes,
380:                 beam_weights=beam_weights,
381:                 codebooks=codebooks,
382:                 scales=scales,
383:                 input_group_index=input_group_index,
384:                 codebook_index=codebook_index,
385:                 best_losses=best_losses,
386:                 best_indices=best_indices,

```

```

aq.py          Fri Mar 22 16:01:31 2024          7
387:             beam_size=beam_size,
388:         )
389:
390:         if verbose:
391:             progressbar.update()
392:             if (input_group_index * num_codebooks + codebook_index) % verbose != 0:
393:                 continue # if update is an integer, compute metrics every (this many) bea
m search steps
394:             best_loss = beam_losses.min(0).values.sum().item() / out_features
395:             info = f"in_group {input_group_index} / {num_in_groups} "
396:             info += f"| codebook {codebook_index} / {num_codebooks} "
397:             if code_penalties is None:
398:                 info += f"| loss {best_loss:.10f}"
399:             else: # un-regularize to restore MSE loss, report sparsity rate
400:                 codebook_ids = torch.arange(num_codebooks, device=beam_losses.device).view
(1, 1, 1, -1)
401:                 best_cand_regularizer = code_penalties[codebook_ids, beam_codes[0]].sum()
/ num_out_groups
402:                 del codebook_ids
403:                 best_loss = best_loss - best_cand_regularizer # report loss without the
regularizer part
404:                 info += f"| loss {best_loss:.5f} | reg {best_cand_regularizer:.5f} |"
405:                 del best_loss
406:                 progressbar.desc = info
407:             return beam_codes[0]
408:
409:
410: @maybe_script
411: def _dequantize_weight(
412:     codes: torch.Tensor, codebooks: torch.Tensor, scales: Optional[torch.Tensor] = None
413: ) -> torch.Tensor:
414:     """
415:     Decode float weights from quantization codes. Differentiable.
416:     :param codes: tensor of integer quantization codes, shape [*dims, num_out_groups, num_in_g
roups, num_codebooks]
417:     :param codebooks: tensor of vectors for each quantization code, [num_codebooks, codebook_s
ize, out_group_size, in_group_size]
418:     :param scales: weight will be multiplied by this factor, must be broadcastble with [*dims,
out_groups, num_in_groups, out_group_size, in_group_size]
419:     :return: reconstructed weight tensor of shape [*dims, num_in_groups*group_size]
420:     """
421:     num_out_groups, num_in_groups, num_codebooks = codes.shape[-3:]
422:     num_codebooks, codebook_size, out_group_size, in_group_size = codebooks.shape
423:     out_features = num_out_groups * out_group_size
424:     in_features = num_in_groups * in_group_size
425:     codebook_offsets = torch.arange(
426:         0, num_codebooks * codebook_size, codebook_size, device=codes.device
427:     ) # shape: [num_codebooks]
428:     reconstructed_weight_flat = F.embedding_bag(
429:         codes.flatten(0, -2) + codebook_offsets, codebooks.flatten(0, 1).flatten(-2, -1), mode
="sum"
430:     ) # [prod(dims) * num_out_groups * num_in_groups, out_group_size * in_group_size]
431:
432:     reconstructed_weight_groupwise = reconstructed_weight_flat.view(
433:         list(codes.shape[:-3]) + [num_out_groups, num_in_groups, out_group_size, in_group_size
]
434:     )
435:     if scales is not None:
436:         reconstructed_weight_groupwise = reconstructed_weight_groupwise.mul(scales)
437:     return reconstructed_weight_groupwise.swapaxes(-3, -2).reshape(list(codes.shape[:-3]) + [o
ut_features, in_features])
438:
439:
440: @maybe_script
441: def _beam_search_squared_errors(
442:     TX: torch.Tensor,
443:     reference_weight: torch.Tensor,
444:     codebooks: torch.Tensor,
445:     scales: Optional[torch.Tensor],
446:     beam_losses: torch.Tensor,
447:     beam_codes: torch.Tensor,
448:     beam_weights: torch.Tensor,
449:     input_group_index: int,
450:     codebook_index: int,
451:     k_best: int,

```

```

452:     code_penalties: Optional[torch.Tensor] = None,
453: ) -> tuple[torch.Tensor, torch.Tensor]:
454:     """
455:     Compute MSE or sum-of-squared-error losses for all possible ways to replace quantization c
odes for one input group
456:     and one codebook. Works in parallel for all output-dimension groups.
457:
458:     :param XTX: pairwise products of input features matmul(X.transpose(), X), shape: [in_featu
res, in_features]
459:     :note: if both XTX *and* beam_losses are divided by dataset size, this function will return
mean squared error
460:     :param reference_weight: original weight matrix that is being quantized, shape: [out_featu
res, in_features]
461:     :param codebooks: look-up tables of codes, shape: [num_codebooks, codebook_size, out_group
_size, in_group_size]
462:     :param scales: weight will be multiplied by this factor, [num_out_groups, num_in_groups, 1
, 1]
463:
464:     :param beam_losses: sum-of-squared-error for each hypothesis in beam and for each output c
hannel;
465:         shape: [beam_size, num_out_groups]
466:     :param beam_codes: a tensor with best weight codes, shape: [beam_size, num_out_groups, num
_in_groups, num_codebooks]
467:     :param beam_weights: a tensor with de-quantized beam_codes, shape: [beam_size, out_feature
s, in_features]
468:     :param input_group_index: an index of one group of in_features that is being re-encoded
469:     :param codebook_index: an index of one codebook for that group of features that is being r
e-encoded
470:     :return: tuple(Tensor, Tensor) of 3d tensor of shape = [beam_size, k_best, num_out_groups]
.
471:         First one is float tensor of losses of k_best lowest square errors for each beam and o
ut_group
472:         Second one is int64 tensor of indices of k_best lowest square errors for each beam and
out_group
473:
474:     :note: The code computes MSE using the square-of-difference expansion
475:         
$$||X@W.T - \sum_i X@(Bi@Ci).T||^2 = ||X@W.T||^2 - 2 \langle X@W.T, \sum_i X@(Bi@Ci).T \rangle + ||\sum_i X@Bi@Ci||^2$$

476:         where X[nsamples,in_features] is calibration data, W[out_features, in_features] is the ref
erence weight,
477:         C[num_codebooks, codebook_size, in_features] are learned codebooks (Ci has shape [codeb
ook_size, out_features])
478:         B[num_codebooks, out_features, codebook_size] are one-hot encoded indices (quantization
codes)
479:         The formula above uses a single group per output "neuron" and a single group.
480:         The algorithm below generalizes the formula for multiple groups and codebooks.
481:
482:         Furthermore, the algorithm does not compute the entire formula. Instead, it begins from so
me baseline loss
483:         and computes the change in loss from changing a single code to every possible altearnative
code.
484:         When computing the changed loss, the algorithm only computes the few affected parts of the
loss formula above.
485:     """
486:     num_codebooks, codebook_size, out_group_size, in_group_size = codebooks.shape
487:     beam_size, num_out_groups, num_in_groups, num_codebooks = beam_codes.shape
488:     out_features = num_out_groups * out_group_size
489:
490:     input_group_slice = slice(input_group_index * in_group_size, (input_group_index + 1) * in_
group_size)
491:
492:     prev_codes_part = beam_codes[:, :, input_group_index, codebook_index] # [beam_size, num_o
ut_groups]
493:
494:     if scales is not None:
495:         scales_part = scales[:, input_group_index % scales.shape[1], :, :] # [num_out_groups,
1, 1]
496:     else:
497:         scales_part = torch.empty(0, device=XTX.device)
498:     prev_part_dequantized = F.embedding(prev_codes_part, codebooks[codebook_index].flatten(-2,
-1)).view(
499:         beam_size, out_features, in_group_size
500:     ) # previous codes de-quantized
501:
502:     prev_weight_part = prev_part_dequantized

```



```

503:     if scales is not None:
504:         prev_weight_part = (
505:             prev_weight_part.view(beam_size, num_out_groups, out_group_size, in_group_size)
506:             .mul(scales_part)
507:             .view(beam_size, out_features, in_group_size)
508:         )
509:
510:     cand_weights = codebooks[codebook_index] # [codebook_size, out_group_size, in_group_size]
, all replacement codes
511:
512:     delta_weight_without_part = reference_weight - beam_weights
513:     delta_weight_without_part[:, :, input_group_slice] += prev_weight_part
514:
515:     # dWTXTX is equivalent to < X @ (W - \sum BiCi except current codebook), X @ SOMETHING >
516:     dWTXTXg = delta_weight_without_part @ XTX[..., input_group_slice] # [beam_size, out_featu
res, in_group_size]
517:     # below: use torch.matmul to compute broadcasted batch matrix multiplication; see matmul d
ocs
518:
519:     XnewBkC_norms_sq = torch.bmm(
520:         (cand_weights.flatten(0, 1) @ XTX[input_group_slice, input_group_slice]).view(
521:             codebook_size, 1, out_group_size * in_group_size
522:         ),
523:         cand_weights.view(codebook_size, out_group_size * in_group_size, 1),
524:     ).reshape(
525:         codebook_size, 1
526:     ) # [codebook_size, num_out_groups]
527:     if scales is not None:
528:         XnewBkC_norms_sq = XnewBkC_norms_sq.mul(scales_part.square()).reshape(1, num_out_groups
))
529:
530:     best_losses = torch.empty(
531:         (beam_size, k_best, num_out_groups), dtype=XTX.dtype, device=XTX.device
532:     ) # shape: [beam_size, k_best, num_out_groups]
533:     best_indices = torch.empty(
534:         (beam_size, k_best, num_out_groups),
535:         dtype=torch.int64,
536:         device=XTX.device,
537:     )
538:     for beam_id in range(beam_size):
539:         dot_products = (
540:             torch.einsum(
541:                 "mg,og->mo",
542:                 cand_weights.reshape(codebook_size, out_group_size * in_group_size),
543:                 dWTXTXg[beam_id].view(num_out_groups, out_group_size * in_group_size),
544:             )
545:             .sub_(
546:                 torch.einsum(
547:                     "og,og->o",
548:                     prev_part_dequantized[beam_id].reshape(num_out_groups, out_group_size * in
_group_size),
549:                     dWTXTXg[beam_id].view(num_out_groups, out_group_size * in_group_size),
550:                 ).view(1, num_out_groups)
551:             )
552:             .view(codebook_size, num_out_groups)
553:         )
554:         if scales is not None:
555:             dot_products = dot_products.mul_(scales_part.reshape(1, num_out_groups))
556:
557:         XoldBkC_norms_sq = torch.bmm(
558:             (prev_weight_part[beam_id] @ XTX[input_group_slice, input_group_slice]).view(
559:                 num_out_groups, 1, out_group_size * in_group_size
560:             ),
561:             prev_weight_part[beam_id].view(num_out_groups, out_group_size * in_group_size, 1),
562:         ).reshape(1, num_out_groups)
563:
564:         # finally, combine them to get MSE
565:         candidate_squared_errors = (
566:             beam_losses[beam_id, None, :] - 2 * dot_products + XnewBkC_norms_sq - XoldBkC_norm
s_sq
567:         ) # shape: [codebook_size, num_out_groups]
568:
569:         if code_penalties is not None:
570:             prev_code_penalties = code_penalties[codebook_index][prev_codes_part[beam_id]] #
[codebook_size]

```

```

571:         candidate_squared_errors[:, :] -= prev_code_penalties[None, :] # refund penalty f
or the replaced code
572:         candidate_squared_errors[:, :] += code_penalties[codebook_index, :, None] # add p
enalty for new code
573:
574:         best_beam_squared_errors, best_beam_indices = torch.topk(
575:             candidate_squared_errors, k_best, dim=0, largest=False, sorted=False
576:         )
577:         best_losses[beam_id] = best_beam_squared_errors
578:         best_indices[beam_id] = best_beam_indices
579:
580:     return best_losses, best_indices
581:
582:
583: @maybe_script
584: def _beam_search_select_best(
585:     beam_codes: torch.Tensor,
586:     beam_weights: torch.Tensor,
587:     codebooks: torch.Tensor,
588:     scales: Optional[torch.Tensor],
589:     input_group_index: int,
590:     codebook_index: int,
591:     best_losses: torch.Tensor,
592:     best_indices: torch.Tensor,
593:     beam_size: int,
594: ) -> Tuple[torch.Tensor, torch.Tensor, torch.Tensor]:
595:     """
596:     Select top-:beam_size: and reorder beam accordingly, return new beam
597:     :param beam_codes: a tensor with best weight codes, shape: [beam_size, num_out_groups, num
_in_groups, num_codebooks]
598:     :param beam_weights: a tensor with de-quantized beam_codes, shape: [beam_size, out_feature
s, in_features]
599:     :param codebooks: a tensor with look-up tables of codes, shape: [num_codebooks, codebook_s
ize, out_group_size, in_group_size]
600:     :param scales: weight will be multiplied by this factor, [num_out_groups, num_in_groups, 1
, 1]
601:
602:     :param input_group_index: an index of one group of in_features that is being re-encoded
603:     :param codebook_index: an index of one codebook for that group of features that is being r
e-encoded
604:     :param best_losses: a 3d tensor of losses of k_best lowest square errors for each beam and
out group,
605:         shape = [beam_size, k_best, num_out_groups]
606:     :param best_indices: a 3d tensor of indices of k_best lowest square errors for each beam a
nd out group,
607:         shape = [beam_size, k_best, num_out_groups]
608:     :param beam_size: how many top hypotheses should be selected
609:
610:     :returns: new (beam_codes, beam_weights, beam_losses)
611:     """
612:     dtype = best_losses.dtype
613:     device = best_losses.device
614:     _prev_beam_size, k_best, num_out_groups = best_losses.shape
615:     _prev_beam_size, out_features, in_features = beam_weights.shape
616:     _prev_beam_size, num_out_groups, num_in_groups, num_codebooks = beam_codes.shape
617:     flat_best = best_losses.flatten(0, 1).topk(dim=0, k=beam_size, largest=False)
618:     best_hypo_source_ids = flat_best.indices // k_best
619:     arange_out_groups = torch.arange(num_out_groups, device=device)
620:     best_hypo_codes = best_indices.flatten(0, 1)[flat_best.indices, arange_out_groups].reshape
(
621:         beam_size, num_out_groups
622:     )
623:     # ^-- shape: [beam_size, num_out_groups]
624:
625:     # reorder beam codes and weights
626:     new_beam_codes = torch.full(
627:         size=(len(best_hypo_codes), num_out_groups, num_in_groups, num_codebooks),
628:         fill_value=-1,
629:         dtype=beam_codes.dtype,
630:         device=device,
631:     ) # [beam_size, num_out_groups, num_in_groups, num_codebooks]
632:     new_beam_weights = torch.empty(len(best_hypo_codes), out_features, in_features, dtype=typ
e, device=device)
633:
634:     for beam_index in range(len(best_hypo_codes)):

```

```

635:         new_beam_codes[beam_index, :, ...] = beam_codes[best_hypo_source_ids[beam_index, :], a
range_out_groups, ...]
636:         new_beam_codes[beam_index, :, input_group_index, codebook_index] = best_hypo_codes[bea
m_index, :]
637:         new_beam_weights[beam_index, :, :] = _dequantize_weight(new_beam_codes[beam_index, ...
], codebooks, scales)
638:
639:         # Note: the code above can be further accelerated by 1) vectorizing loop and ...
640:         # ... 2) updating new_beam_weights only for the chosen input group
641:         return new_beam_codes, new_beam_weights, flat_best.values
642:
643:
644: @maybe_script
645: def _channelwise_squared_error(XTX: torch.Tensor, weight: torch.Tensor, reference_weight: torch
h.Tensor):
646:     """
647:     Compute per-channel squared error between X @ weight_or_weights and X @ reference_weight
648:     :param XTX: pairwise products of input features matmul(X.transpose(), X), shape: [in_featu
res, in_features]
649:     :note: if XTX is divided by dataset size, this function will return *mean* squared error
650:     :param weight: predicted/reconstructed weights of shape [*dims, out_features, in_features]
651:     :param reference_weight: reference weight of shape [out_features, in_features]
652:     :return: per-channel squared errors of shape [*dims, out_features]
653:     """
654:     XW_norm_square = torch.matmul(weight[..., :, None, :], (weight @ XTX)[..., :, :, None]).fl
atten(-3)
655:     XWreference_norm_square = torch.bmm(reference_weight[:, None, :], (reference_weight @ XTX)
[:, :, None]).flatten(-3)
656:     dot_product = torch.matmul((reference_weight @ XTX)[:, None, :], weight[..., :, :, None]).
flatten(-3)
657:     return XW_norm_square - 2 * dot_product + XWreference_norm_square
658:
659:
660: @torch.no_grad()
661: def init_aq_kmeans(
662:     reference_weight: torch.Tensor,
663:     *,
664:     num_codebooks: int,
665:     out_group_size: int,
666:     in_group_size: int,
667:     codebook_size: int,
668:     verbose: bool = False,
669:     use_faiss: bool = False,
670:     max_points_per_centroid: Optional[int] = None,
671:     max_iter: int = 1000,
672:     devices: Optional[List[torch.device]] = None,
673:     **kwargs,
674: ):
675:     """
676:     Create initial codes and codebooks using residual K-means clustering of weights
677:     :params reference_weight, num_codebooks, out_group_size, in_group_size, nbits, verbose: sa
me as in QuantizedWeight
678:     :params use_faiss whether to use faiss implementation of kmeans or pure torch
679:     :params max_point_per_centroid maximum data point per cluster
680:     :param kwargs: any additional params are forwarded to fit_kmeans
681:     """
682:     out_features, in_features = reference_weight.shape
683:     num_out_groups = out_features // out_group_size
684:     num_in_groups = in_features // in_group_size
685:     weight_residue = (
686:         reference_weight.reshape(num_out_groups, out_group_size, num_in_groups, in_group_size)
687:         .clone()
688:         .swapaxes(-3, -2)
689:         .reshape(num_out_groups * num_in_groups, out_group_size * in_group_size)
690:     )
691:     codebooks = []
692:     codes = []
693:
694:     if max_points_per_centroid is not None:
695:         print("Clustering:", max_points_per_centroid * codebook_size, "points from", weight_re
sidue.shape[0])
696:
697:     for _ in trange(num_codebooks, desc="initializing with kmeans") if verbose else range(num_
codebooks):
698:         if use_faiss:

```

```
699:         codebook_i, codes_i, reconstructed_weight_i = fit_faiss_kmeans(
700:             weight_residue,
701:             k=codebook_size,
702:             max_iter=max_iter,
703:             gpu=(weight_residue.device.type == "cuda"),
704:             max_points_per_centroid=max_points_per_centroid,
705:         )
706:     else:
707:         chosen_ids = None
708:         if max_points_per_centroid is not None:
709:             chosen_ids = torch.randperm(weight_residue.shape[0], device=weight_residue.dev
ice)[
710:                 : max_points_per_centroid * codebook_size
711:             ]
712:         codebook_i, _, _ = fit_kmeans(
713:             weight_residue if chosen_ids is None else weight_residue[chosen_ids, :],
714:             k=codebook_size,
715:             max_iter=max_iter,
716:             devices=devices,
717:             **kwargs,
718:         )
719:         codes_i, reconstructed_weight_i = find_nearest_cluster(weight_residue, codebook_i,
devices=devices)
720:
721:         codes_i = codes_i.reshape(num_out_groups, num_in_groups, 1)
722:         codebook_i = codebook_i.reshape(1, codebook_size, out_group_size, in_group_size)
723:         weight_residue -= reconstructed_weight_i
724:         codes.append(codes_i)
725:         codebooks.append(codebook_i)
726:         del reconstructed_weight_i
727:         codebooks = torch.cat(codebooks, dim=0)
728:         codes = torch.cat(codes, dim=-1)
729:         return codes, codebooks
```