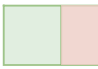
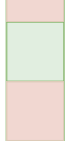


def masked_softmax(X, valid_lens): #@save

def _sequence_mask(
 X,
 valid_len,
 value=0
):
 maxlen = X.size(1)

maxlen = length_k

mask = torch.arange(
 (maxlen),
 dtype=torch.float32,
 device=X.device
)[None, :] < valid_len[:, None]

($B * h * \text{length}_q, \text{length}_k$)(16, 2)
mask =  ($1, 2$) <  ($16, 1$)
torch.arange(maxlen)[None, :] valid_lens[:, None]
($1, \text{length}_k$) ($B * h * \text{length}_q, 1$)

X[~mask] = value

let ~mask position pad by 0

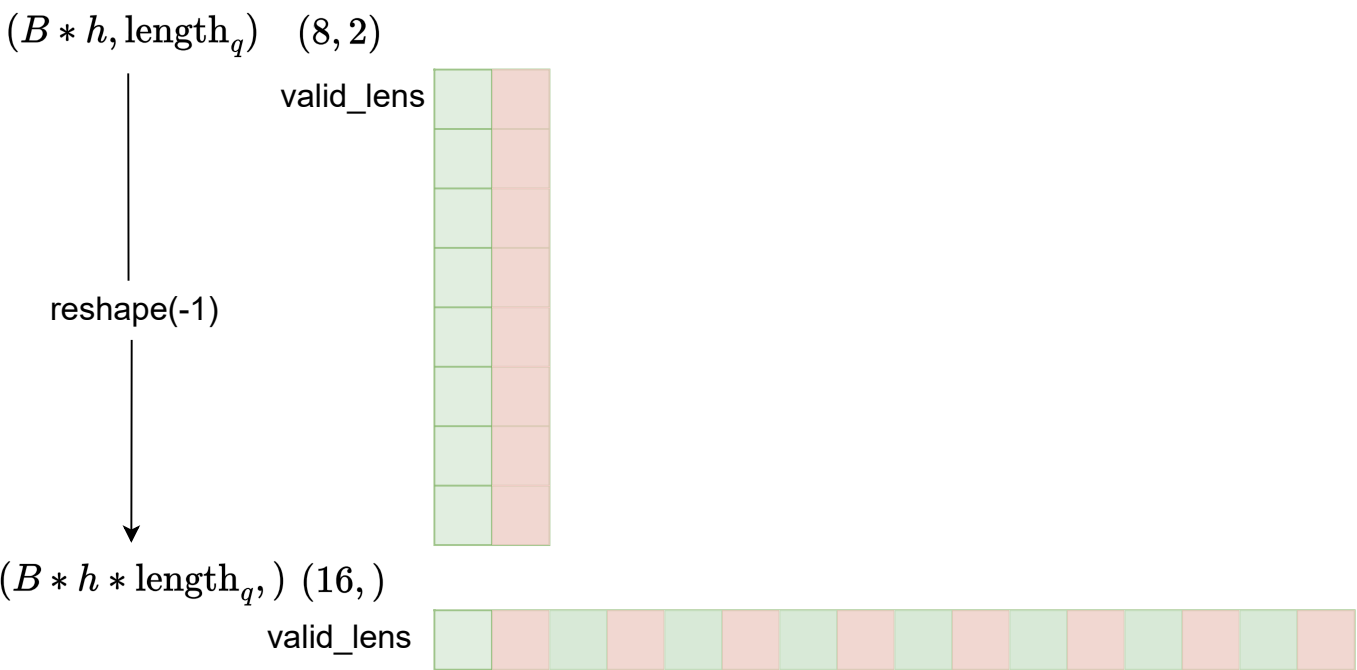
return X

($B * h * \text{length}_q, \text{length}_k$)(16, 2)

if valid_lens is None:
 return nn.functional.softmax(X, dim=-1)
else:
 shape = X.shape

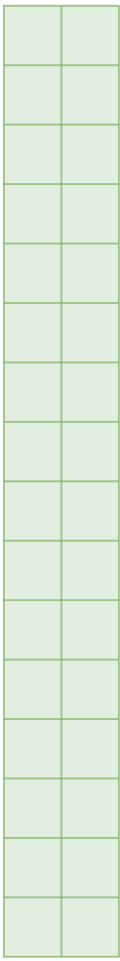
shape = ($B * h, \text{length}_q, \text{length}_k$) (8, 2, 2)

if valid_lens.dim() == 1:
 valid_lens = torch.repeat_interleave(
 valid_lens, shape[1]
)
else:
 valid_lens = valid_lens.reshape(-1)



X = _sequence_mask(
 X.reshape(-1, shape[-1]),
 valid_lens,
 value=-1e6
)

($B * h * \text{length}_q, \text{length}_k$)(16, 2)
 X



return nn.functional.softmax(
 X.reshape(shape), dim=-1
)

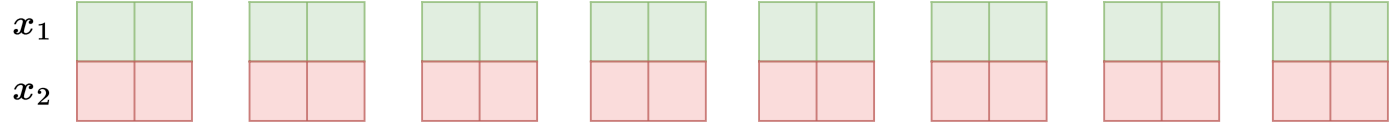
($B * h, \text{length}_q, \text{length}_k$) (8, 2, 2)
 X

softmax to the last dim (length_k)

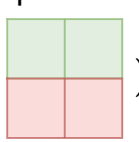
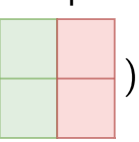
class DotProductAttention(nn.Module): #@save
 def __init__(self, dropout):
 super().__init__()
 self.dropout = nn.Dropout(dropout)

def forward(self, queries, keys, values, valid_lens=None):
 queries=keys=values=X $d = d_{model}/h$
 d = queries.shape[-1]

($B * h, \text{length}, d_{model}/h$) (8, 2, 2)
 X



scores = torch.bmm(
 queries,
 keys.transpose(1, 2)
) / math.sqrt(d)


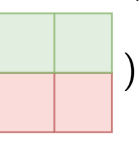
($B * h, \text{length}_q, \text{length}_k$)(8, 2, 2)
queries keys.transpose(1,2)
(8, ) × (8, )
scores = $\frac{\quad}{\sqrt{d}}$



self.attention_weights = masked_softmax(
 scores,
 valid_lens
)

the dim is the same as scores,
we do the mask at length_k dim

return torch.bmm(
 self.dropout(self.attention_weights),
 values
)

($B * h, \text{length}_q, \text{length}_k$) ($B * h, \text{length}_k, d_{model}/h$)
(8, ) × (8, )
out = $\frac{\quad}{\sqrt{d}}$

($B * h, \text{length}_q, d_{model}/h$)(8, 2, 2)