


```

out: softmax output for each data point, shape (batch, n_class)
axis: the attention value on channels

```

```

TODO
1. [Given] pass each channel of x through the corresponding beat_net, then
   We will discard the attention (alpha and beta) outputs for now
   Using ModuleList for self.beat_nets/rhythm_nets is necessary for the
2. [Given] stack the output from 1 together into a tensor of shape (batch,
3. pass result from 2 and k_freq through attention module, to get the agg
   You might need to do use k_freq.permute() to tweak the shape of k_freq
4. pass aggregated result from 3 through the final fully connected layer.
5. Apply Softmax to normalize output to a probability distribution (over 2)

"""
def forward(self, x, k_beats, k_rhythms, k_freq):
    new_x = [None for _ in range(self.n_channels)]

    # 1. pass each channel of x through the corresponding beat_net, then rhythm_net
    # Do not collect the attention outputs
    # We used ModuleList so that the gradient would propagate

    for i in range(self.n_channels):
        tx, _ = self.beat_nets[i](x[i], k_beats[i])
        tx, _ = self.rhythm_nets[i](tx, k_rhythms[i])

    # 2. stack the output into a tensor with shape (batch, n_channels, rhythm_out
    # output shape: [128,8]
    # stack along dim = 1 to get [128,4,8]
    x = torch.stack(new_x, 1)

    # CODE

    # 3. pass result from 2 and k_freq through attention module, to get the aggregated result
    # k_freq.shape: torch.Size([4, 17, 1]) => 17, 4, 1
    # out: shape of (-1, input_features)

    out, gamma = self.attn(x, torch.permute(k_freq, (1, 0, 2)))

    # 4. pass aggregated result from 3 through the final fully connected layer.
    out = self.fc(out)

    # 5. Apply Softmax to normalize output to a probability distribution (over 2)
    out = torch.softmax(out, dim = 1)

    # END

    return out, gamma

```

```

In [17]:
"""
AUTOGRADE CELL. DO NOT MODIFY THIS.
"""
_B, _M, _T = 17, 59, 109
testm = FreqNet(_B*_M*_T, _M*_T, _T*_T)
assert isinstance(testm.attn, KnowledgeAttn), "Should use one KnowledgeAttn Module"
assert isinstance(testm.fc, nn.Linear) and testm.fc.weight.shape == torch.Size([2,8])
assert isinstance(testm.beat_nets, nn.ModuleList), "beat_nets has to be a ModuleList"

_out, gamma = testm(torch.randn(4, _B, _M*_T), torch.randn(4, _B, _M*_T), torch.randn(4, _B, 4, 1))
assert _out.shape == torch.Size([_B, 4, 1]), "The attention's dimension is incorrect"
del _testm, _out, _gamma, _B, _M, _T

```

```

In [ ]:

```

3 Training and Evaluation [15 points]

In this part we will define the training procedures, train the model, and evaluate the model on the test set.

train_model parameters:

- model: The instance of FreqNet that we are training
 - train_dataloader: the DataLoader of the training data
 - n_epoch: number of epochs to train
 - lr: learning rate
 - device: cpu or gpu/cuda
- return/output:**
- _model: trained model
 - _loss_history: recorded training loss history - should be just a list of float
- train_model tasks:**
- I. Specify the optimizer (*optimizer*) to be optim.Adam
- II. Specify the loss function (*loss_func*) to be CrossEntropyLoss
- III. Within the loop, do the normal training procedures:
- A. pass the input through the model
 - B. pass the output through loss_func to compute the loss
 - C. zero out currently accumulated gradient, use loss.backward to backprop the gradients, then call optimizer.step

eval_model tasks:

- returns:**
- **pred_all:** prediction of model on the dataloader.
 - Should be an 2D numpy float array where the second dimension has length 2.
 - **Y_test:** truth labels. Should be an numpy array of ints
- tasks:**
1. evaluate the model using on the data in the dataloader.
 2. Add all the prediction and truth to the corresponding list
 3. Convert pred_all and Y_test to numpy arrays (of shape (n_data_points, 2))

ADAM

params (Iterable) – iterable of parameters to optimize or dicts defining parameter groups

```

torch.optim.Adam(
    params,
    lr=0.001,
    betas=(0.9, 0.999),
    eps=1e-08,
    weight_decay=0,
    amsgrad=False, *,
    foreach=None,
    maximize=False,
    capturable=False,
    differentiable=False,
    fused=None
)

```

CrossEntropyLoss

```

torch.nn.CrossEntropyLoss(
    weight=None,
    ignore_index=- 100,
    reduction='mean',
    label_smoothing=0.0
)

```

loss between input logits and target

input: unnormalized logits for each class.

input has to be a Tensor of Size (C), for unbached input (minibatch, C) or (minibatch, c, d1, d2, ... dk) for d = 1 to dim.

target:

Class indices in the range (0, C), C = number of classes (*preferred*)

or

Probabilities for each class

```

In [18]:
def train_model(model, train_dataloader, n_epoch=5, lr=0.003, device=None):
    import torch.optim as optim

    device = device or torch.device('cpu')
    model.train()

    loss_history = []

    # CODE

    # I. Specify the optimizer ("optimizer") to be optim.Adam

    optimizer = optim.Adam(model.parameters(), lr=lr)
    # II. Specify the loss function (loss_func) to be CrossEntropyLoss
    loss_func = nn.CrossEntropyLoss()

    # END

    for epoch in range(n_epoch):
        curr_epoch_loss = []
        for (X, K_beat, K_rhythm, K_freq), Y in train_dataloader:
            print("##### X.shape: #####", X.shape)
            # CODE

            # III. Within the loop, do the normal training procedures:
            pred = model(X, K_beat, K_rhythm, K_freq)[0] # Y = model(X) ## ??? # A.
            loss = loss_func(pred, Y.long()) # B. pass the output through loss_func to

            # C.
            optimizer.zero_grad() # i. zero out currently accumulated gradient,
            loss.backward() # ii. use loss.backward to backprop the gradient
            optimizer.step() # iii. call optimizer.step

            # END

            curr_epoch_loss.append(loss_func(pred, Y).data.numpy())
        print(f"Epoch{epoch}: curr_epoch_loss={np.mean(curr_epoch_loss)}")
        loss_history += curr_epoch_loss
    return model, loss_history

```

```

In [ ]:

```

```

In [19]:
device = torch.device('cpu')
n_epoch = 4
lr = 0.003
n_channel = 4
n_dim=3000
T=50

model = FreqNet(n_channel, n_dim, T)
model = model.to(device)

model, loss_history = train_model(model, train_loader, n_epoch=n_epoch, lr=lr, device=device)
pred, truth = eval_model(model, test_loader, device=device)
sgd_to_pickle((pred, truth), "delivariable.pkl")

```

```

In [20]:
pred, truth = eval_model(model, test_loader, device=device)

```

```

In [24]:
pred.shape, truth.shape

```

```

In [21]:
%debug

```

```

In [26]:
f1_score(y_true, y_pred)

```

```

In [27]:
def evaluate_predictions(truth, pred):
    """
    TODO: Evaluate the performance of the predictoin via AUROC, and F1 score
    each prediction in pred is a vector representing (p,0, p,1).
    When defining the scores we are interested in detecting class 1 only, ie 0, 1
    (Hint: use roc_auc_score and f1_score from sklearn.metrics, be sure to read their
    return: auroc, f1
    """
    from sklearn.metrics import roc_auc_score, f1_score

    auroc = roc_auc_score(truth, pred[:, 1])
    pred = np.argmax(pred, axis=1)
    f1_f1_score(truth, pred)

    return auroc, f1

```

```

In [28]:
"""
AUTOGRADE CELL. DO NOT MODIFY THIS.
"""
pred, truth = eval_model(model, test_loader, device=device)
auroc, f1 = evaluate_predictions(truth, pred)
print(f"AUROC={auroc} and F1={f1}")

assert auroc > 0.8 and f1 > 0.7, "Performance is too low {}. Something's probably off

```

```

In [ ]:

```

```

In [ ]:

```