


```
In [ ]: class Attention(nn.Module):
    def __init__(self, feature_dim, step_dim, bias=True, **kwargs):
        super(Attention, self).__init__(**kwargs)

        self.supports_masking = True

        self.bias = bias
        self.feature_dim = feature_dim
        self.step_dim = step_dim
        self.steps_dim = 0

        weight = torch.zeros(feature_dim, 1)
        nn.init.xavier_uniform_(weight)
        self.weight = nn.Parameter(weight)

        if bias:
            self.b = nn.Parameter(torch.zeros(step_dim))

    def forward(self, x, mask=None):
        feature_dim = self.feature_dim
        step_dim = self.step_dim

        eij = torch.mm(
            x.contiguous().view(-1, feature_dim),
            self.weight
        ).view(-1, step_dim)

        if self.bias:
            eij = eij + self.b

        eij = torch.tanh(eij)
        a = torch.exp(eij)

        if mask is not None:
            a = a * mask

        a = a / torch.sum(a, 1, keepdim=True) + 1e-10

        weighted_sum = x * torch.unsqueeze(a, -1)
        return torch.sum(weighted_input, 1)

class NeuralNet(nn.Module):
    def __init__(self):
        super(NeuralNet, self).__init__()

        fc_layer1 = 16
        fc_layer1 = 16

        self.embedding = nn.Embedding(max_features, embed_size)
        self.embedding.weight = nn.Parameter(torch.tensor(embedding_matrix, dtype=torch.float32))
        self.embedding.requires_grad = False

        self.embedding_dropout = nn.Dropout2d(0.1)
        self.lstm = nn.LSTM(embed_size, hidden_size, bidirectional=True, batch_first=True)
        self.gru = nn.GRU(hidden_size * 2, hidden_size, bidirectional=True, batch_first=True)

        self.lstm2 = nn.LSTM(hidden_size * 2, hidden_size, bidirectional=True, batch_first=True)

        self.lstm_attention = Attention(hidden_size * 2, maxlen)
        self.gru_attention = Attention(hidden_size * 2, maxlen)

        self.linear = nn.Linear(hidden_size*8+3, fc_layer1) #643:80 - 483:60 - 323:40
        self.relu = nn.ReLU()
        self.dropout = nn.Dropout(0.1)
        self.fc = nn.Linear(fc_layer1**2, fc_layer1)
        self.out = nn.Linear(fc_layer1, 1)
        self.lincaps = nn.Linear(Num_capsule * Dim_capsule, 1)
        self.caps_layer = Capsule_layer()

    def forward(self, x):

        Capsule(num_capsule=10, dim_capsule=10, routings=4, share_weights=True)(x)

        h_embedding = self.embedding(x[0])
        h_embedding = torch.squeeze(
            self.embedding_dropout(torch.unsqueeze(h_embedding, 0)))

        h_lstm, _ = self.lstm(h_embedding)
        h_gru, _ = self.gru(h_lstm)

        #Capsule Layer
        content3 = self.caps_layer(h_gru)
        content3 = self.dropout(content3)
        batch_size = content3.size(0)
        content3 = content3.view(batch_size, -1)
        content3 = self.relu(self.lincaps(content3))

        #Attention Layer
        h_lstm_atten = self.lstm_attention(h_lstm)
        h_gru_atten = self.gru_attention(h_gru)

        # global average pooling
        avg_pool = torch.mean(h_gru, 1)
        # global max pooling
        max_pool, _ = torch.max(h_gru, 1)

        f = torch.tensor(x[1], dtype=torch.float).cuda()

        # [512,160]
        conc = torch.cat((h_lstm_atten, h_gru_atten, content3, avg_pool, max_pool, f), 1)
        conc = self.relu(self.linear(conc))
        conc = self.dropout(conc)
        out = self.out(conc)

        return out
```

Training

The method for training is borrowed from <https://www.kaggle.com/hengsheng/bytorch-starter>

```
In [ ]: class MyDataset(Dataset):
    def __init__(self, dataset):
        self.dataset = dataset

    def __getitem__(self, index):
        data, target = self.dataset[index]

        return data, target, index

    def __len__(self):
        return len(self.dataset)

In [ ]: def sigmoid(x):
    return 1 / (1 + np.exp(-x))

# matrix for the out-of-fold predictions
train_preds = np.zeros((len(x_train)))
# matrix for the predictions on the test set
test_preds = np.zeros((len(df_test)))

# always call this before training for deterministic results
seed_everything()

x_test_cuda = torch.tensor(x_test, dtype=torch.long).cuda()
# test_f = torch.utils.data.TensorDataset(x_test_cuda, f)
# test_loader = torch.utils.data.DataLoader(test, batch_size=batch_size, shuffle=False)

x_test_cuda = torch.tensor(x_test, dtype=torch.long).cuda()
test = torch.utils.data.TensorDataset(x_test_cuda)
test_loader = torch.utils.data.DataLoader(test, batch_size=batch_size, shuffle=False)

avg_losses_f = []
avg_val_losses_f = []

In [ ]: for i, (train_idx, valid_idx) in enumerate(splits):
    # split data in train / validation according to the KFold indexes
    # also, convert them to a torch tensor and store them on the GPU (done with .cuda())
    x_train = np.array(y_train)
    y_train = np.array(y_train)
    features = np.array(features)

    x_train_fold = torch.tensor(x_train[train_idx.astype(int)], dtype=torch.long).cuda()
    y_train_fold = torch.tensor(y_train[train_idx.astype(int)], np.newaxis), dtype=torch.float32).cuda()

    kfold_y_valid_features = features[train_idx.astype(int)]
    kfold_x_valid_features = features[valid_idx.astype(int)]
    x_val_fold = torch.tensor(x_train[valid_idx.astype(int)], dtype=torch.long).cuda()
    y_val_fold = torch.tensor(y_train[valid_idx.astype(int)], np.newaxis), dtype=torch.float32).cuda()

    # model = BiLSTM(lstm_layer=2, hidden_dim=40, dropout=DROPOUT).cuda()
    model = NeuralNet()

    # make sure everything in the model is running on the GPU
    model.cuda()

    # define binary cross entropy loss
    # note that the model returns logit to take advantage of the log-sum-exp trick
    # for numerical stability in the loss
    loss_fn = torch.nn.BCEWithLogitsLoss(reduction='sum')

    step_size = 300
    base_lr, max_lr = 0.001, 0.003
    optimizer = torch.optim.Adam(filter(lambda p: p.requires_grad, model.parameters()),
                                     lr=max_lr)

    #####
    scheduler = Cyclical(optimizer, base_lr=base_lr, max_lr=max_lr,
                          step_size=step_size, mode='exp_range',
                          gamma=0.99994)
    #####

    train = torch.utils.data.TensorDataset(x_train_fold, y_train_fold)
    valid = torch.utils.data.TensorDataset(x_val_fold, y_val_fold)

    train = MyDataset(train)
    valid = MyDataset(valid)

    #Who need to shuffle the data again here. Shuffling happens when splitting for k-folds.
    train_loader = torch.utils.data.DataLoader(train, batch_size=batch_size, shuffle=True)
    valid_loader = torch.utils.data.DataLoader(valid, batch_size=batch_size, shuffle=False)

    print(f'Fold {i + 1}')
    for epoch in range(n_epochs):
        # set train mode of the model. This enables operations which are only applied during training
        # like dropout
        start_time = time.time()
        model.train()

        avg_loss = 0
        for i, (x_batch, y_batch, index) in enumerate(train_loader):
            # Forward pass: compute predicted y by passing x to the model:
            #####
            f = kfold_x_valid_features[index]
            y_pred = model(x_batch, f)
            #####
            #####

            if scheduler:
                scheduler.batch_step()
            #####

            # Compute and print loss.
            loss = loss_fn(y_pred, y_batch)

            # Before the backward pass, use the optimizer object to zero all of the
            # gradients for the Tensors it will update (which are the learnable weights
            # of the model)
            optimizer.zero_grad()

            # Backward pass: compute gradient of the loss with respect to model parameters
            loss.backward()

            # Calling the step function on an Optimizer makes an update to its parameters
            optimizer.step()
            avg_loss += loss.item() / len(train_loader)

        # set evaluation mode of the model. This disabled operations which are only applied during training
        # like dropout
        model.eval()

        # predict all the samples in y_val_fold batch per batch
        valid_preds_fold = np.zeros((x_val_fold.size(0)))
        test_preds_fold = np.zeros((len(df_test)))

        avg_val_loss = 0
        for i, (x_batch, y_batch, index) in enumerate(valid_loader):
            f = kfold_x_valid_features[index]
            y_pred = model(x_batch, f).detach()

            avg_val_loss += loss_fn(y_pred, y_batch).item() / len(valid_loader)
            valid_preds_fold[i * batch_size:(i+1) * batch_size] = sigmoid(y_pred.cpu().numpy())[:, 0]

        elapsed_time = time.time() - start_time
        print(f'Epoch {i/(n_epochs)} \t loss={:.4f} \t val_loss={:.4f} \t time={:.2f}s'.format(
            epoch + 1, n_epochs, avg_loss, avg_val_loss, elapsed_time))
        avg_losses_f.append(avg_loss)
        avg_val_losses_f.append(avg_val_loss)
        # predict all samples in the test set batch per batch
        for i, (x_batch,) in enumerate(test_loader):
            f = test_features[i * batch_size:(i+1) * batch_size]
            y_pred = model(x_batch, f).detach()

            test_preds_fold[i * batch_size:(i+1) * batch_size] = sigmoid(y_pred.cpu().numpy())[:, 0]

        train_preds[valid_idx] = valid_preds_fold
        test_preds += test_preds_fold / len(splits)

    print(f'All \t loss={:.4f} \t val_loss={:.4f} \t '.format(np.average(avg_losses_f), np.average(avg_val_losses_f)))

# x_train, x_test_f, y_train, y_test_f
```

Find final Threshold

Borrowed from: <https://www.kaggle.com/ziliwang/baseline-pytorch-bilstm>

```
In [ ]: def bestThreshold(y_train, train_preds):
    tmp = [0, 0]
    for idx, cur, max in tqdm(np.arange(0.1, 0.501, 0.01)):
        tmp[0] = f1_score(y_train, np.array(train_preds)>tmp[0])
        if tmp[0] > tmp[2]:
            delta = tmp[0]
            tmp[2] = tmp[0]
    print('best threshold is {:.4f} with F1 score: {:.4f}'.format(delta, tmp[2]))
    return delta
delta = bestThreshold(y_train, train_preds)

In [ ]: submission = df_test[['qid']].copy()
submission['prediction'] = (test_preds > delta).astype(int)
submission.to_csv('submission.csv', index=False)

In [ ]: !head submission.csv
```