P.S. не забудьте, что если вы работаете в Colab, то вы можете поменять среду выполнения на GPU/TPU! Вопросы по заданию/материалам: @Nestyme In [1]: !pip3 install timit-utils==0.9.0 !pip3 install torchaudio ! wget https://ndownloader.figshare.com/files/10256148 !unzip -q 10256148 In [2]: import timit\_utils as tu import os import librosa import numpy as np from tqdm import tqdm import torch import torch.nn as nn from torch.optim import Adam import torch.nn.functional as F import matplotlib.pyplot as plt from sklearn.metrics import accuracy\_score import IPython \_TIMIT\_PATH = 'data/lisa/data/timit/raw/TIMIT' Задание 1 Загрузите данные для обучения. Для этого: 1. Скачайте датасет TIMIT (см семинар) 2. Соберите пары "голос" — "класс возраста" также, как на семинаре собирались пары "голос" — "пол". Аудиодорожки сконвертируйте в мелспектрограммы при помощи torchaudio либо librosa P.S. вы можете использовать свою реализацию, а можете предложенную (см следующие ячейки) In [3]: import timit\_utils as tu import os import librosa import numpy as np from tqdm import tqdm import torch as t class timit\_dataloader: def \_\_init\_\_(self, data\_path=\_TIMIT\_PATH, train\_mode=True, age\_mode=True): self.doc\_file\_path = os.path.join(data\_path, 'DOC', 'SPKRINFO.TXT') self.corpus = tu.Corpus(data\_path) with open(self.doc\_file\_path) as f: self.id\_age\_dict = dict( [(tmp.split(' ')[0], 86 - int(tmp.split(' ')[5].split('/')[-1].replace('??', '50'))) \ for tmp in f.readlines()[39:]]) if train\_mode: self.trainset = self.create\_dataset('train', age\_mode=age\_mode) self.validset = self.create\_dataset('valid', age\_mode=age\_mode) self.testset = self.create\_dataset('test', age\_mode=age\_mode) def return\_age(self, id): return self.id\_age\_dict[id] def return\_data(self): return self.trainset, self.validset, self.testset def return\_test(self): return self.testset def create\_dataset(self, mode, age\_mode=False): global people assert mode in ['train', 'valid', 'test'] if mode == 'train': people = [self.corpus.train.person\_by\_index(i) for i in range(350)] if mode == 'valid': people = [self.corpus.train.person\_by\_index(i) for i in range(350, 400)] if mode == 'test': people = [self.corpus.test.person\_by\_index(i) for i in range(150)] spectrograms\_and\_targets = [] for person in tqdm(people): try: target = self.return\_age(person.name) for i in range(len(person.sentences)): spectrograms\_and\_targets.append( self.preprocess\_sample(person.sentence\_by\_index(i).raw\_audio, target, age\_mode=True)) except: print(person.name, target) X, y = map(np.stack, zip(\*spectrograms\_and\_targets)) X = X.transpose([0, 2, 1]) # to [batch, time, channels] return X, y @staticmethod def spec\_to\_image(spec, eps=1e-6): mean = spec.mean() std = spec.std() spec\_norm = (spec - mean) / (std + eps) spec\_min, spec\_max = spec\_norm.min(), spec\_norm.max() spec\_scaled = 255 \* (spec\_norm - spec\_min) / (spec\_max - spec\_min) spec\_scaled = spec\_scaled.astype(np.uint8) return spec\_scaled @staticmethod def clasterize\_by\_age(age): **if** age <=25: return 0 **if** 25 < age <= 40: return 1 **if** age > 40: return 2 def preprocess\_sample(self, amplitudes, target, age\_mode=False, sr=16000, max\_length=150): spectrogram = librosa.feature.melspectrogram(amplitudes, sr=sr, n\_mels=128, fmin=1, fmax=8192)[:, :max\_length] spectrogram = np.pad(spectrogram, [[0, 0], [0, max(0, max\_length - spectrogram.shape[1])]], mode='constant') target = self.clasterize\_by\_age(target) return self.spec\_to\_image(np.float32(spectrogram)), target def preprocess\_sample\_inference(self, amplitudes, sr=16000, max\_length=150, device='cpu'): spectrogram = librosa.feature.melspectrogram(amplitudes, sr=sr, n\_mels=128, fmin=1, fmax=8192)[:, :max\_length] spectrogram = np.pad(spectrogram, [[0, 0], [0, max(0, max\_length - spectrogram.shape[1])]], mode='constant') spectrogram = np.array([self.spec\_to\_image(np.float32(spectrogram))]).transpose([0, 2, 1]) return t.tensor(spectrogram, dtype=t.float).to(device, non\_blocking=True) class dataloader: def \_\_init\_\_(self, spectrograms, targets): self.data = list(zip(spectrograms, targets)) def next\_batch(self, batch\_size, device): indices = np.random.randint(len(self.data), size=batch\_size) input = [self.data[i] for i in indices] source = [line[0] for line in input] target = [line[1] for line in input] return self.torch\_batch(source, target, device) @staticmethod def torch\_batch(source, target, device): return tuple( t.tensor(val, dtype=t.float).to(device, non\_blocking=True) for val in [source, target] @staticmethod def padd\_sequences(lines, pad\_token=0): lengths = [len(line) for line in lines] max\_length = max(lengths) return np.array( line + [pad\_token] \* (max\_length - lengths[i]) for i, line in enumerate(lines) Простая сверточная сеть, ее можно дотюнить или поменять по желанию In [4]: import torch import torch.nn as nn import torch.nn.functional as F class Model(nn.Module): def \_\_init\_\_(self, window\_sizes=(3, 4, 5)): super(Model, self).\_\_init\_\_() self.convs = nn.ModuleList([ nn.Conv2d(1, 128, [window\_size, 128], padding=(window\_size - 1, 0)) for window\_size in window\_sizes self.bn = nn.BatchNorm1d(128) self.fc = nn.Linear(128 \* len(window\_sizes), 3) self.drop = nn.Dropout(p=0.15) def forward(self, x): x = torch.unsqueeze(x, 1) # [B, C, T, E] Add a channel dim.xs = []for conv in self.convs: x2 = F.relu(conv(x)) # [B, F, T, 1]x2 = torch.squeeze(x2, -1) # [B, F, T] $x2 = F.max_{pool1d}(x2, x2.size(2)) # [B, F, 1]$ xs.append(x2)x = torch.cat(xs, 2) # [B, F, window]x = self.bn(x)x = self.drop(x)# FC x = x.view(x.size(0), -1) # [B, F \* window]logits = self.fc(x) # [B, class] probs = torch.sigmoid(logits) #.view(-1) return probs In [5]: \_timit\_dataloader = timit\_dataloader() train, valid, test = \_timit\_dataloader.return\_data() trainset = dataloader(\*train) validset = dataloader(\*valid) testset = dataloader(\*test) 100% 350/350 [00:38<00: 00, 9.08it/s] 100% 50/50 [00:05<00: 00, 8.54it/s] 100% 150/150 [00:17<00: 00, 8.64it/s] Задание 2 1. Обучите свой классификатор категории возраста 2. Попробуйте улучшить результат. Можно попробовать усложнить сетку, подвигать границы категорий, поискать новые данные, что угодно, кроме учиться на тесте :) 3. Какой подход оказался самым эффективным? Как думаете, почему? 4. Как считаете, где можно было бы применить такой классификатор в качестве вспомогательной задачи? In [6]: device = torch.device('cuda' if torch.cuda.is\_available() else 'cpu') print(f'using {device} mode') patience = 500 best\_loss = 1000 cnt = 0 using cuda mode In [7]: import gc gc.collect() torch.cuda.empty\_cache() torch.cuda.reset\_max\_memory\_allocated() torch.cuda.synchronize() C:\Users\BIT\.conda\envs\deep\lib\site-packages\torch\cuda\memory.py:231: FutureWarning: torch.cuda.reset\_max\_memory\_allocated now calls torch.cuda.reset\_peak\_memory\_stats, which resets /all/ peak memory stats. warnings.warn( In [8]: model = Model() if device == torch.device('cuda'): model.cuda() model.cpu() model.train() Out[8]: Model( (convs): ModuleList( (0): Conv2d(1, 128, kernel\_size=[3, 128], stride=(1, 1), padding=(2, 0)) (1): Conv2d(1, 128, kernel\_size=[4, 128], stride=(1, 1), padding=(3, 0)) (2): Conv2d(1, 128, kernel\_size=[5, 128], stride=(1, 1), padding=(4, 0)) (bn): BatchNorm1d(128, eps=1e-05, momentum=0.1, affine=True, track\_running\_stats=True) (fc): Linear(in\_features=384, out\_features=3, bias=True) (drop): Dropout(p=0.15, inplace=False) In [9]: BATCH\_SIZE = 64 optimizer = Adam( [p for p in model.parameters() if p.requires\_grad], betas=(0.9, 0.999), lr=0.005, eps=1e-5 def model\_loss(probs, targets): return nn.CrossEntropyLoss(weight=1-torch.Tensor([0.24,0.63,0.13]).cuda())(probs, targets) #return nn.CrossEntropyLoss()(probs, targets) In [10]: from IPython.display import clear\_output import math min\_loss = np.inf min\_acc = np.inf train\_loss = [] train\_history = [] train\_acc = [] valid\_acc = [] valid\_history = [] model.train() for i in range(1500): optimizer.zero\_grad() input, target = trainset.next\_batch(BATCH\_SIZE, device=device) out = model(input) loss = model\_loss(out, target.long()) loss.backward() optimizer.step() **if** i % 50 == 0: model.eval() with torch.no\_grad(): optimizer.zero\_grad() input, target = validset.next\_batch(BATCH\_SIZE, device=device) out = model(input) valid\_loss = model\_loss(out, target.long()) out, target = out.cpu().detach().numpy(), target.cpu().detach().numpy() out = np.argmax(out,axis=1) val\_acc = accuracy\_score(target,out) valid\_acc.append(val\_acc) valid\_history.append(valid\_loss.item()) print(f'accuracy\_score:{val\_acc}') print("i {}, valid {}".format(i, valid\_loss.item())) print("\_\_\_\_\_") model.train() fig, ax = plt.subplots(nrows=1, ncols=3, figsize=(16, 6)) clear\_output(True) if train\_history is not None: ax[0].plot(train\_history, label='train history') ax[0].set\_xlabel('Epoch') ax[0].set\_title('Train loss') ax[0].legend() if valid\_history is not None: ax[1].plot(valid\_history, label='valid history') ax[1].legend() ax[1].set\_xlabel('Epoch') ax[1].set\_title('Valid loss') if valid\_acc is not None: ax[2].plot(valid\_acc, label='valid accuracy') ax[2].legend() ax[2].set\_xlabel('Epoch') ax[2].set\_title('Valid accuracy') plt.legend() plt.show() train\_history.append(loss.item()) if i % 50 == 0 and best\_loss > valid\_loss.item(): best\_loss = valid\_loss.item() cnt = 0 else: cnt += 1 if cnt > patience: break print('training finished') Valid loss Valid accuracy Train loss 1.20 — train history — valid accuracy 0.65 1.15 -0.60 0.55 1.10 0.50 1.05 0.45 0.40 0.8 1.00 0.35 0.7 0.30 valid history 12 14 100 200 300 400 500 600 700 2 10 12 14 Epoch training finished In [11]: model.eval() for i in range(30): with torch.no\_grad(): optimizer.zero\_grad() input, target = testset.next\_batch(BATCH\_SIZE, device=device) out = model(input) test\_loss = model\_loss(out, target.long()) out, target = out.cpu().detach().numpy(), target.cpu().detach().numpy() out = np.argmax(out,axis=1) print(f'accuracy\_score:{accuracy\_score(out, target)}') print("\_\_\_\_\_") print('testing finished') accuracy\_score:0.515625 accuracy\_score:0.34375 accuracy\_score:0.4375 accuracy\_score:0.515625 accuracy\_score:0.296875 accuracy\_score:0.34375 accuracy\_score:0.5 accuracy\_score:0.4375 accuracy\_score:0.375 accuracy\_score:0.40625 accuracy\_score:0.5 accuracy\_score:0.59375 accuracy\_score:0.46875 accuracy\_score:0.3125 accuracy\_score:0.40625 accuracy\_score:0.40625 accuracy\_score:0.40625 accuracy\_score:0.359375 accuracy\_score:0.421875 accuracy\_score:0.453125 accuracy\_score:0.453125 accuracy\_score:0.4375 accuracy\_score:0.484375 accuracy\_score:0.453125 accuracy\_score:0.578125 accuracy\_score:0.484375 accuracy\_score:0.46875 accuracy\_score:0.34375 accuracy\_score:0.515625 accuracy\_score:0.5 testing finished Выводы: 1) Постарался обучить свой классификатор возраста по голосу, разделив людей на группы (до 25, от 25 до 40, после 40). 2) Усложнил семинарскую сеть со свёртками с различными kernel\_size (3,4,5) и оптимизатором Adam с Ir=0.005, добавив BatchNorm, Dropout с p=0.15 и линейный слой в конец модели для классификации на 3 класса; взял CE loss с применением весов (0.24,0.63,0.13), из-за выхода модели в 3 класса и уменьшения негативного влияния дисбаланса классов. Количество итераций взял 1500 с предварительным выходом и batch size = 64. Результат получился немного лучше константного предсказания, по большей части, как мне кажется, из-за примитивности предобработки и слабой модели (касательно использования различных идей). Также пытался двигать границы для уменьшения дисбаланса, но из-за малых различий в голосах модель не могла обучиться; пытался использовать Factorized-TDNN, добавив линейный слой в конец, но сеть плохо обучалась. 3) Добавление BatchNorm'а оказался самым эффективным подходом, то есть результат сильно улучшился и проблем с переобучением на самый

распространённый класс стало меньше, из-за того что линейная модель (линейный слой) лучше работает с нормированными данными; также

возможно неплохое улучшение результата могут дать новые данные и последующий подбор границ с наиболее различающимися голосами.

4) В голосовом помощнике, например, чтобы модель понимала пользуется ли ей ребёнок, молодой человек или старик, с помощью данной

(например, компания-маркетплейс).

информации можно отфильтровывать ответы, которые обычно являются неподходящими для данной категории, и таким образом делать более

с автоматом, определить какие её продукты могли бы быть ему предположительно нужны, после чего сообщать данную информацию оператору

хорошие рекомендации. Также с помощью данного классификатора какая-либо компания, могла бы используя голос человека, пока он разговаривает

**Spoken Language Processing** 

В этом задании предлагается обучить классификатор возраста по голосу (пример с тем, как это можно сделать для пола см. в семинаре)

Подумайте, как лучше предсказывать возраст (может быть разбить на группы?) и какой лосс использовать