

캡스톤 디자인 ‘딥페이크 탐지’

#8. Adversarial training

김지수, 김민지, 민지민

지난 캡스톤 회의 내용 -블랙박스 공격 조사-

denoise, sharpening, gaussian noise, salt and pepper, jpeg 으로 만들어진 데이터셋에 대한 강도별 성능 측정

=> 이번주 캡스톤 내용

성능이 잘 나오지 않은 sharpening, gaussian noise, salt and pepper 를 이용해서 adversarial training 진행

데이콘 1등 코드 오류난 부분 코랩으로 실행

→ 또 다른 import 오류로 실행 불가

Adversarial training

```
1 model = xception(num_out_classes=2, dropout=0.5)
2 # model = models.xception()
3 print("=> creating model '{}'.format('xception'))
4 # model = model.cuda(args.gpu)
5
6 # fn = 'deepfake_c0_xception.pkl'
7 # assert os.path.isfile(fn), 'wrong path'
8
9 # model.load_state_dict(torch.load(fn))
10 best_model = torch.load("/content/drive/My Drive/Anti-forensics/deepfake_c0_xception_tuned.pth.tar")
11 model.load_state_dict(best_model['state_dict'])
12 print("=> model weight best_model is loaded")
13
14 model.to(device)
15
16 criterion = nn.CrossEntropyLoss().cuda()
17 optimizer = optim.Adam(model.parameters(), lr=0.001, betas=(0.9, 0.999), eps=1e-08)
18 scheduler = lr_scheduler.StepLR(optimizer, step_size=5, gamma=0.5)
```

저번주에 학습한 모델 가져와
adversarial training 진행

	salt_pepper (strong)	salt_pepper (medium)	salt_pepper (weak)	sharpening (strong)	sharpening (medium)	sharpening (weak)
gaussian noise strong	loss:5.4151 <u>accuracy:0.246</u>	loss: 6.9820 <u>accuracy:0.245</u>	loss:4.5988 <u>accuracy:0.263</u>	loss: accuracy:	loss: accuracy:	loss: accuracy:
gaussian noise med	loss: 8.1304 <u>accuracy:0.291</u>	loss: 10.0488 <u>accuracy:0.291</u>	loss : 6.2867 <u>accuracy :0.297</u>	loss: 0.0779 <u>accuracy:0.968</u>	loss: 0.2596 <u>accuracy:0.901</u>	loss: 0.5444 <u>accuracy:0.809</u>
gaussian noise weak	loss: 0.0615 <u>accuracy:0.988</u>	loss:1.2993 <u>accuracy:0.542</u>	loss: 3.0358 <u>accuracy:0.450</u>	loss: 0.0358 <u>accuracy:0.991</u>	loss: 0.1095 <u>accuracy:0.963</u>	loss: 0.4228 <u>accuracy:0.881</u>

	sharpening (weak)	sharpening (medium)	sharpening (strong)	gaussian noise (weak)	gaussian noise (medium)	gaussian noise (strong)
salt & pepper noise (weak)	loss - 8.2891 <u>acc - 0.255</u>	loss - 9.0158 <u>acc - 0.251</u>	loss -10.5005 <u>acc - 0.250</u>	loss - 11.8990 <u>acc - 0.250</u>	loss - 12.4867 <u>acc - 0.250</u>	loss - 2.7501 <u>acc - 0.442</u>
salt & pepper noise (medium)	loss - 0.2123 <u>acc - 0.912</u>	loss - 0.1200 <u>acc - 0.949</u>	loss - 0.1163 <u>acc - 0.950</u>	loss - 0.2564 <u>acc - 0.890</u>	loss - 0.1162 <u>acc - 0.950</u>	loss - 0.1163 <u>acc - 0.950</u>
salt & pepper noise (strong)	loss - 0.0899 <u>acc - 0.964</u>	loss - 0.0289 <u>acc - 0.989</u>	loss - 0.0280 <u>acc - 0.989</u>	loss - 5.1495 <u>acc - 0.312</u>	loss - 1.5534 <u>acc - 0.380</u>	loss - 0.0276 <u>acc - 0.989</u>

	salt & pepper (weak)	salt & pepper (medium)	salt & pepper (strong)	gaussian noise (weak)	gaussian noise (medium)	gaussian noise (strong)
sharpening (weak)	<u>acc - 0.005</u>	<u>acc - 0.017</u>	<u>acc - 0.135</u>	<u>acc - 0.677</u>	<u>acc - 0.936</u>	<u>acc - 1.0</u>
sharpening (medium)	<u>acc - 0</u>	<u>acc - 0</u>	<u>acc - 0.010</u>	<u>acc - 0.011</u>	<u>acc - 0.057</u>	<u>acc - 1.0</u>
sharpening (strong)	<u>acc - 0.069</u>	<u>acc - 0</u>	<u>acc - 0</u>	<u>acc - 0.506</u>	<u>acc - 0.625</u>	<u>acc - 1.0</u>