

1 Задание «Предсказание карт внимания», часть 2

1.1 Описание

Предлагается реализовать нейросетевой алгоритм построения карт внимания

Требования: * Модель должна быть реализована на фреймворке PyTorch * Для обучения разрешается использовать только обучающую выборку, внешние данные использовать нельзя. Однако не запрещено (и, наоборот, приветствуется) использовать аугментации (размножение датасета) и transfer-learning (использование предобученной модели) с дообучением * Обработка одного кадра должна быть возможна на видеокарте с объемом памяти 24Gb * Соблюдать [кодекс чести](#). Виновные будут найдены и наказаны

1.1.1 Оценивание

1. Каждый участник может представить не более 1 алгоритма для финального тестирования
2. Тестирование будет проводиться на закрытой тестовой выборке, содержащей N ($N < 20$) тестовых видео. В каждом > 200 кадров.
3. В качестве метрик будут использованы:
 - Normalized Scanpath Saliency (NSS)
 - Similarity score (SIM)
 - Pearson's Correlation Coefficient (CC)

[Подробнее про метрики](#)

4. По итогам тестирования будет составлена общая таблица результатов по каждой из метрик
5. Место алгоритма определяется по формуле:

$$\text{Place}_{\text{algo}} = \frac{\text{Place}_{\text{NSS}} + \text{Place}_{\text{SIM}} + \text{Place}_{\text{CC}}}{3}$$

6. Баллы участника зависят от места его алгоритма:

$$\text{Score}_{\text{stud}} = \text{score}(\text{Place}_{\text{algo}})$$

GPU 0: NVIDIA A100-SXM4-80GB (UUID: GPU-39f56421-0a72-e5a8-8dde-39060aadf095)

1.2 1. Подготовка данных

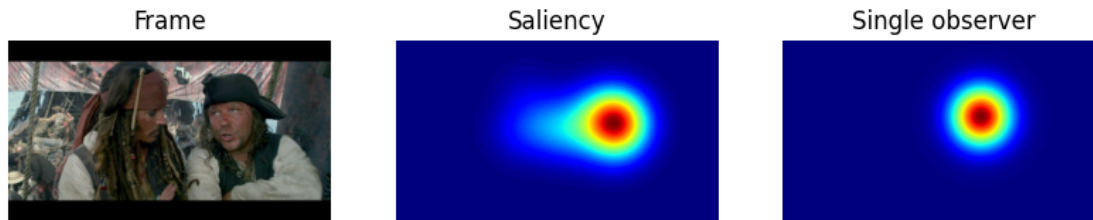
[Дублирование датасета с страницы задания на Google Диск для пользователей Colab \[7.2 GB\]](#)

1.2.1 Датасеты и даталоадеры

Вы можете полностью менять эту реализацию, использовать более одного кадра, добавлять transforms и т.д.!

Посмотрим на случайный пример из датасета

(-0.5, 383.5, 215.5, -0.5)



1.3 2. Создадим класс нашей модели

Так как данных мало, предлагается использовать технику **transfer learning**, используя нейросеть, предобученную на сегментацию изображений, например, **Deeplabv3**.

Модель построим из трех частей: * Сверточный Encoder, например, ResNet50 * Пирамидальный пулинг ASPP состоящий из нескольких параллельных сверток с разным рецептивным полем благодаря dilation convolution * Сверточный Decoder, получающий финальное предсказание

Архитектура Deeplab:

Dilation convolution:

В модели, доступной в библиотеке `torchvision.models`, Decoder получает изображение в низком разрешении и имеет 21 класс (количество выходных каналов):

```
Conv2d(256, 21, kernel_size=(1, 1), stride=(1, 1))
```

В нашей же задаче требуется предсказания одноканального изображения с разрешением, как у исходного. Поэтому сделаем свой декодер, постепенно повышая разрешение в 2 раза. На последнем слое применим нормализацию в 0-1.

Encoder ожидает на вход 3х-канальное изображение:

```
Conv2d(3, 64, kernel_size=(7, 7), stride=(2, 2), padding=(3, 3), bias=False)
```

Документация PyTorch: <https://pytorch.org/docs/stable/index.html>

Ваша модель должна быть описана в файле `saliency_single_observer.py`, этот файл **СДАЕТСЯ** в проверяющую систему!

Файл `saliency_single_observer.py` должен быть самодостаточным, т.е. содержать все необходимые `import`-ы и функции для создания и запуска модели.

Layer (type:depth-idx)	Output Shape	Param #
SingleObserverSaliencyModel	[8, 1, 216, 384]	--
└IntermediateLayerGetter: 1-1	[8, 2048, 27, 48]	--
└Conv2d: 2-1	[8, 64, 108, 192]	12,544
└BatchNorm2d: 2-2	[8, 64, 108, 192]	128
└ReLU: 2-3	[8, 64, 108, 192]	--
└MaxPool2d: 2-4	[8, 64, 54, 96]	--
└Sequential: 2-5	[8, 256, 54, 96]	--
└Bottleneck: 3-1	[8, 256, 54, 96]	75,008

└─Bottleneck: 3-2	[8, 256, 54, 96]	70,400
└─Bottleneck: 3-3	[8, 256, 54, 96]	70,400
└─Sequential: 2-6	[8, 512, 27, 48]	--
└─Bottleneck: 3-4	[8, 512, 27, 48]	379,392
└─Bottleneck: 3-5	[8, 512, 27, 48]	280,064
└─Bottleneck: 3-6	[8, 512, 27, 48]	280,064
└─Bottleneck: 3-7	[8, 512, 27, 48]	280,064
└─Sequential: 2-7	[8, 1024, 27, 48]	--
└─Bottleneck: 3-8	[8, 1024, 27, 48]	1,512,448
└─Bottleneck: 3-9	[8, 1024, 27, 48]	1,117,184
└─Bottleneck: 3-10	[8, 1024, 27, 48]	1,117,184
└─Bottleneck: 3-11	[8, 1024, 27, 48]	1,117,184
└─Bottleneck: 3-12	[8, 1024, 27, 48]	1,117,184
└─Bottleneck: 3-13	[8, 1024, 27, 48]	1,117,184
└─Sequential: 2-8	[8, 2048, 27, 48]	--
└─Bottleneck: 3-14	[8, 2048, 27, 48]	6,039,552
└─Bottleneck: 3-15	[8, 2048, 27, 48]	4,462,592
└─Bottleneck: 3-16	[8, 2048, 27, 48]	4,462,592
─ASPP: 1-2	[8, 256, 27, 48]	--
└─ModuleList: 2-9	--	--
└─Sequential: 3-17	[8, 256, 27, 48]	524,800
└─ASPPConv: 3-18	[8, 256, 27, 48]	4,719,104
└─ASPPConv: 3-19	[8, 256, 27, 48]	4,719,104
└─ASPPConv: 3-20	[8, 256, 27, 48]	4,719,104
└─ASPPPooling: 3-21	[8, 256, 27, 48]	524,800
└─Sequential: 2-10	[8, 256, 27, 48]	--
└─Conv2d: 3-22	[8, 256, 27, 48]	327,680
└─BatchNorm2d: 3-23	[8, 256, 27, 48]	512
└─ReLU: 3-24	[8, 256, 27, 48]	--
└─Identity: 3-25	[8, 256, 27, 48]	--
─Sequential: 1-3	[8, 128, 27, 48]	--
└─Conv2d: 2-11	[8, 128, 27, 48]	294,912
└─BatchNorm2d: 2-12	[8, 128, 27, 48]	256
└─ReLU: 2-13	[8, 128, 27, 48]	--
─Sequential: 1-4	[8, 64, 54, 96]	--
└─Conv2d: 2-14	[8, 64, 54, 96]	73,728
└─BatchNorm2d: 2-15	[8, 64, 54, 96]	128
└─ReLU: 2-16	[8, 64, 54, 96]	--
─Sequential: 1-5	[8, 32, 108, 192]	--
└─Conv2d: 2-17	[8, 32, 108, 192]	18,432
└─BatchNorm2d: 2-18	[8, 32, 108, 192]	64
└─ReLU: 2-19	[8, 32, 108, 192]	--
─Conv2d: 1-6	[8, 1, 216, 384]	289

1.4 3. Зададим функцию потерь

Будем использовать дивергенцию [Кульбака-Лейблера](#), сравнивая предсказанное распределение с эталонной картой внимания. Вы так же можете использовать другие функции потерь при обучении.

1.5 4. Обучим модель

Пайплайн обучения: * Определить **таргет**. В нашей задаче это эталонная карта внимания * Определить **функцию потерь (loss)**. Используем KLD из предыдущего этапа * Выбрать **оптимизатор**. Чтобы всё быстро заработало, возьмем Adam/AdamW с $lr=3e-4$

```
100%|██████████| 50/50 [00:03<00:00, 13.46it/s, loss=0.518]
```

```
100%|██████████| 3/3 [00:00<00:00, 11.13it/s]
```

```
| Epoch: 0 | Val Loss: 0.45534321665763855 | Train Loss: 0.6122541737556457
```

```
100%|██████████| 50/50 [00:03<00:00, 14.24it/s, loss=0.399]
```

```
100%|██████████| 3/3 [00:00<00:00, 12.27it/s]
```

```
| Epoch: 1 | Val Loss: 0.40258293350537616 | Train Loss: 0.4304718714952469
```

```
100%|██████████| 50/50 [00:03<00:00, 14.30it/s, loss=0.416]
```

```
100%|██████████| 3/3 [00:00<00:00, 11.91it/s]
```

```
| Epoch: 2 | Val Loss: 0.288569043080012 | Train Loss: 0.4019764757156372
```

```
100%|██████████| 50/50 [00:03<00:00, 13.85it/s, loss=0.286]
```

```
100%|██████████| 3/3 [00:00<00:00, 12.61it/s]
```

```
| Epoch: 3 | Val Loss: 0.3885070780913035 | Train Loss: 0.41201020300388336
```

```
100%|██████████| 50/50 [00:03<00:00, 13.71it/s, loss=0.509]
```

```
100%|██████████| 3/3 [00:00<00:00, 9.88it/s]
```

```
| Epoch: 4 | Val Loss: 0.41874605417251587 | Train Loss: 0.4122008267045021
```

```
100%|██████████| 50/50 [00:03<00:00, 14.39it/s, loss=0.317]
```

```
100%|██████████| 3/3 [00:00<00:00, 10.68it/s]
```

```
| Epoch: 5 | Val Loss: 0.2235511839389801 | Train Loss: 0.3856837475299835
```

```
100%|██████████| 50/50 [00:03<00:00, 14.37it/s, loss=0.372]
```

```
100%|██████████| 3/3 [00:00<00:00, 9.26it/s]
```

```
| Epoch: 6 | Val Loss: 0.2144584208726883 | Train Loss: 0.3431571587920189
```

```
100%|██████████| 50/50 [00:03<00:00, 13.09it/s, loss=0.387]
100%|██████████| 3/3 [00:00<00:00, 12.93it/s]
```

| Epoch: 7 | Val Loss: 0.3624043216307958 | Train Loss: 0.3641256356239319

```
100%|██████████| 50/50 [00:03<00:00, 14.26it/s, loss=0.338]
100%|██████████| 3/3 [00:00<00:00, 11.96it/s]
```

| Epoch: 8 | Val Loss: 0.2868318359057109 | Train Loss: 0.3584832298755646

```
100%|██████████| 50/50 [00:03<00:00, 13.67it/s, loss=0.344]
100%|██████████| 3/3 [00:00<00:00, 10.26it/s]
```

| Epoch: 9 | Val Loss: 0.5922253131866455 | Train Loss: 0.3423562940955162

1.6 5. Протестируем модель

Сохраняем веса модели (этот файл **СДАЕТСЯ** в проверяющую систему)

В качестве метрик будут использованы: * Normalized Scanpath Saliency (NSS) * Similarity score (SIM) * Pearson's Correlation Coefficient (CC)

[Подробнее про метрики](#)

Помимо описания класса модели в **СДАВАЕМОМ** файле `saliency_single_observer.py` необходимо описать класс `SingleObserverSaliencyEvaluator`, выполняющий логику загрузки модели и инференса на заданном видео с сохранением предсказаний.

Функция принимает путь до входных последовательностей кадров и путь до выходной папки предсказаний.

В этой функции нужно описать процесс загрузки весов модели получение карт внимания для всех входных видеопоследовательностей. Если ваша модель использует более одного кадра или требует дополнительных преобразований входа, реализуйте их внутри этой функции.

```
100%|██████████| 450/450 [00:34<00:00, 13.12it/s]
100%|██████████| 450/450 [00:36<00:00, 12.44it/s]
100%|██████████| 450/450 [00:31<00:00, 14.34it/s]
100%|██████████| 450/450 [00:30<00:00, 14.98it/s]
100%|██████████| 450/450 [00:30<00:00, 14.73it/s]
100%|██████████| 450/450 [00:29<00:00, 15.51it/s]
100%|██████████| 450/450 [00:31<00:00, 14.12it/s]
100%|██████████| 450/450 [00:29<00:00, 15.33it/s]
100%|██████████| 450/450 [00:27<00:00, 16.58it/s]
```

```
100%|██████████| 3/3 [00:00<00:00, 59074.70it/s]
100%|██████████| 450/450 [00:05<00:00, 75.84it/s]
100%|██████████| 450/450 [00:06<00:00, 74.76it/s]
100%|██████████| 450/450 [00:05<00:00, 75.47it/s]
100%|██████████| 450/450 [00:05<00:00, 75.26it/s]
100%|██████████| 450/450 [00:06<00:00, 73.40it/s]
100%|██████████| 450/450 [00:06<00:00, 74.19it/s]
```

```
100%|██████████| 450/450 [00:06<00:00, 74.82it/s]
100%|██████████| 450/450 [00:06<00:00, 73.99it/s]
100%|██████████| 450/450 [00:06<00:00, 74.54it/s]
```

```
{'sim': 0.501303672660819, 'nss': 1.8504584636522106, 'cc': 0.5417532095532906}
```

```
100%|██████████| 90/90 [00:01<00:00, 53.51it/s]
```

1.7 6. Дальнейшие шаги

Для улучшения качества вашей модели можно попробовать следующие этапы: * Сделать более информативную валидацию, например, добавив подсчет тестовых метрик прямо в validation loop * Попробовать обучать большее число эпох, использовать регуляризацию, например, через Dropout * Добавить skip-connection в Encoder-Decoder * Сделать аугментации ([albumentations](#)). Будьте аккуранты с преобразованиями, которые могут потенциально изменить эталонные карты внимания. Начать можно с горизонтальных отражений * Попробовать другие архитектуры, функции потерь и стратегии обучения * Использовать информацию из более чем одного кадра (3D Conv/LSTM/GRU/любой другой способ агрегации). Обратите внимание: даже если ваш метод требует окно из кадров, тестирование всё равно будет учитывать все кадры видео. Во время тестирования вы можете искусственно дублировать первый кадр для накопления нужной ширины окна для предсказания. * Поискать методы, решающие похожие задачи * Пофантазировать и вдохновиться