

# MRO - zadanie 6: Konwolucje jednak niepotrzebne...

Piotr Zawisła

16 Grudzień 2022

## 1 Wykorzystane technologie

Kod pisałem w języku Python z wykorzystaniem frameworka Pytorch.

## 2 Przygotowanie danych

Transformację danych wykonałem używając modułu **torchvision.transform**. Pobrałem dataset, a następnie podmieniłem atrybut transform na złożone transformacje, które zajmują się augmentacją danych.

```
1 cifar_train_ds = datasets.CIFAR10(  
2     root="data",  
3     train=True,  
4     download=True,  
5     transform=Compose([ToTensor(),  
6                         ConvertImageDtype(torch.float)]),  
7     target_transform=Lambda(  
8         lambda y: torch.zeros(10, dtype=torch.float).scatter_(0, torch.tensor(y),  
9         value=1)  
10    )  
11 )  
12 cifar_train_ds.transform = Compose([  
13     ToTensor(),  
14     RandomHorizontalFlip(p=0.5),  
15     RandomResizedCrop(size=32, scale=(0.8, 1.0), ratio=(0.9, 1.1)),  
16     Normalize(  
17         mean=tuple((cifar_train_ds.data / 255).mean(axis=(0, 1, 2))),  
18         std=tuple((cifar_train_ds.data / 255).std(axis=(0, 1, 2)))  
19     ),  
20     ConvertImageDtype(torch.float)  
21 ])
```

Listing 1: Transformacje wejściowe dla zbioru CIFAR10

## 3 Cięcie obrazu na fragmenty

Aby podzielić obraz na patche, dwa razy użyłem metody **unfold**, a następnie zamieniłem wymiary, tak aby uzyskać kształt nadający się do dalszego przetwarzania.

```
1 def get_square_patches_from_batch(X: torch.Tensor, size: int) -> torch.Tensor:  
2     """Returns tensor of patches with shape [B, H // size, W // size, C, size, size]  
3     """  
4     assert X.ndim == 4, 'ndim should be 4'  
5     batch_size, n_channels, n_rows, n_cols = X.shape  
6     step = size  
7     patches = X.unfold(2, size, step).unfold(3, size, step).permute(0, 2, 3, 1, 4, 5)  
8     return patches
```

Listing 2: Funkcja tnąca obraz na kwadratowe patche.



Rysunek 1: Ptak (prawdopodobnie kazuar)



Rysunek 2: Ptak pocięty na patche o rozmiarze 8

## 4 Definicja modelu

Zaprojektowałem model *Vision Transformer*a, jak określone w poleceniu. Niestety, nie udało mi się wyświetlić computational graph sieci, ale raczej działa jak należy :).

```
1 class ToPatches(nn.Module):
2     """
3     Transforms image to flattened patches.
4
5     [B, C, H, W] -> [B, H // size, W // size, C, size, size] ->
6     [B, N (n_patches), P (size_of_flattened_patch)]
7     """
8
9     def __init__(self, patch_size=4) -> None:
10         super().__init__()
11         self.patch_size: int = patch_size
12
13     def forward(self, X: torch.Tensor) -> torch.Tensor:
14         patches = get_square_patches_from_batch(X, size=self.patch_size)
15         flat_patches = patches.flatten(start_dim=3).flatten(start_dim=1, end_dim=2)
16         return flat_patches
```

Listing 3: Definicja modułu zwracającego pocięte i spłaszczone patche.

```
1 class VisionTransformer(nn.Module):
2     def __init__(self, emb_size=256, n_heads=8, dropout_rate=0.2) -> None:
3         super().__init__()
4
5         self.emb_size: int = emb_size
6         self.n_heads: int = n_heads
7         self.dropout_rate: float = dropout_rate
8
9         self.layer_norm_1 = nn.LayerNorm(normalized_shape=self.emb_size)
10        self.mh_attention = nn.MultiheadAttention(
11            embed_dim=self.emb_size, num_heads=self.n_heads, batch_first=True
12        )
13
14        self.layer_norm_2 = nn.LayerNorm(normalized_shape=self.emb_size)
15        self.mlp_1 = nn.Linear(in_features=self.emb_size, out_features=2*self.emb_size
16    )
17        self.gelu = nn.GELU()
18        self.dropout_1 = nn.Dropout(self.dropout_rate)
19
20        self.mlp_2 = nn.Linear(in_features=2*self.emb_size, out_features=self.emb_size
21    )
22        self.dropout_2 = nn.Dropout(self.dropout_rate)
23
24    def forward(self, X: torch.Tensor) -> torch.Tensor:
25        X_norm = self.layer_norm_1(X)
26
27        mh_attention_output, _ = self.mh_attention(
28            query=X_norm, key=X_norm, value=X_norm
29        )
30
31        residual_add_1 = mh_attention_output + X
32
33        X_mlp_1 = self.layer_norm_2(residual_add_1)
34        X_mlp_1 = self.mlp_1(X_mlp_1)
35        X_mlp_1 = self.gelu(X_mlp_1)
36        result_mlp_1 = self.dropout_1(X_mlp_1)
37
38        X_mlp_2 = self.mlp_2(result_mlp_1)
39        result_mlp_2 = self.dropout_2(X_mlp_2)
40
41        residual_add_2 = residual_add_1 + result_mlp_2
42
43        return residual_add_2
```

Listing 4: Definicja modułu pojedynczego transformera z multi-head attention.

```

1 class CifarViTClassifier(nn.Module):
2     def __init__(
3         self, n_classes=10, patch_size=4,
4         n_attention_heads=8, n_transformers=6,
5         emb_size=256, image_size=32, dropout_rate=0.2
6     ) -> None:
7
8         super().__init__()
9
10        self.n_classes: int = n_classes
11        self.patch_size: int = patch_size
12        self.n_attention_heads: int = n_attention_heads
13        self.n_transformers: int = n_transformers
14        self.emb_size: int = emb_size
15        self.image_size: int = image_size
16        self.dropout_rate: float = dropout_rate
17
18        self.patch_dim: int = 3 * self.patch_size**2
19        self.n_patches: int = (self.image_size // self.patch_size)**2
20
21        # learnable parameters
22        self.pos_embedding = nn.Parameter(torch.randn(1, self.n_patches + 1, self.
emb_size))
23        self.cls_token = nn.Parameter(torch.randn(1, 1, self.emb_size))
24
25        # image -> patches -> patches_embedding
26        self.to_patch_embedding = nn.Sequential(
27            ToPatches(patch_size=self.patch_size),
28            nn.Linear(in_features=self.patch_dim, out_features=self.emb_size)
29        )
30
31        self.input_dropout = nn.Dropout(p=self.dropout_rate)
32
33        # self.n_transformers x transformers
34        self.transformers_seq = nn.Sequential(*[
35            VisionTransformer(
36                emb_size=self.emb_size,
37                n_heads=self.n_attention_heads,
38                dropout_rate=self.dropout_rate
39            ) for _ in range(self.n_transformers)
40        ])
41
42        # linear head
43        self.mlp_head = nn.Sequential(
44            nn.LayerNorm(self.emb_size),
45            nn.Linear(self.emb_size, self.n_classes)
46        )
47
48
49    def forward(self, X: torch.Tensor) -> torch.Tensor:
50        batch_size = X.shape[0]
51        embedded_patches = self.to_patch_embedding(X)
52
53
54        cls_tokens = self.cls_token.expand(batch_size, -1, -1)
55        embedded_patches_with_cls = torch.cat((cls_tokens, embedded_patches), dim=1)
56        embedded_patches_with_cls += self.pos_embedding
57
58        trans_X = self.input_dropout(embedded_patches_with_cls)
59        trans_output = self.transformers_seq(trans_X)
60
61        cls_tokens_out = trans_output[:, 0, :]
62
63        result = self.mlp_head(cls_tokens_out)
64        return result

```

Listing 5: Definicja modelu ViT.

## 5 Trening

Wytrenowałem model przez 160 epok, z optimizerem **AdamW** oraz z schedulerem **MultiStepLR**. Trening trwał około 2,6 godziny z użyciem karty GTX 1660 Ti.

```
1 def train_step(  
2     model: nn.Module,  
3     data_loader: DataLoader,  
4     loss_fn: nn.CrossEntropyLoss,  
5     optimizer  
6 ) -> None:  
7  
8     num_batches = len(data_loader)  
9     epoch_loss = 0.  
10    epoch_correct_preds = 0  
11    model.train()  
12    t_start = perf_counter()  
13  
14    for batch, (X, y) in enumerate(data_loader):  
15        X, y = X.to('cuda'), y.to('cuda')  
16        optimizer.zero_grad()  
17        prediction = model(X)  
18        loss = loss_fn(prediction, y)  
19        loss.backward()  
20        optimizer.step()  
21        epoch_loss += loss.item()  
22        correct_preds = (prediction.argmax(dim=1) == y.argmax(dim=1)).count_nonzero().  
23        item()  
24        epoch_correct_preds += correct_preds  
25  
26    t_end = perf_counter()  
27    epoch_train_time = t_end - t_start  
28    epoch_avg_loss = epoch_loss / num_batches  
29    epoch_avg_acc = epoch_correct_preds / (num_batches * data_loader.batch_size)  
30  
31    return epoch_avg_loss, epoch_avg_acc, epoch_train_time
```

Listing 6: Funkcja wykonująca jedną epokę treningu.

```
1 def val_step(  
2     model: nn.Module,  
3     data_loader: DataLoader,  
4     loss_fn: nn.CrossEntropyLoss  
5 ) -> None:  
6  
7     num_batches = len(data_loader)  
8     epoch_loss = 0.  
9     epoch_correct_preds = 0  
10    model.eval()  
11    with torch.no_grad():  
12        for batch, (X, y) in enumerate(data_loader):  
13            X, y = X.to('cuda'), y.to('cuda')  
14            prediction = model(X)  
15            loss = loss_fn(prediction, y)  
16            epoch_loss += loss.item()  
17            correct_preds = (prediction.argmax(dim=1) == y.argmax(dim=1)).  
18            count_nonzero().item()  
19            epoch_correct_preds += correct_preds  
20  
21    epoch_avg_loss = epoch_loss / num_batches  
22    epoch_avg_acc = epoch_correct_preds / (num_batches * data_loader.batch_size)  
23  
24    return epoch_avg_loss, epoch_avg_acc
```

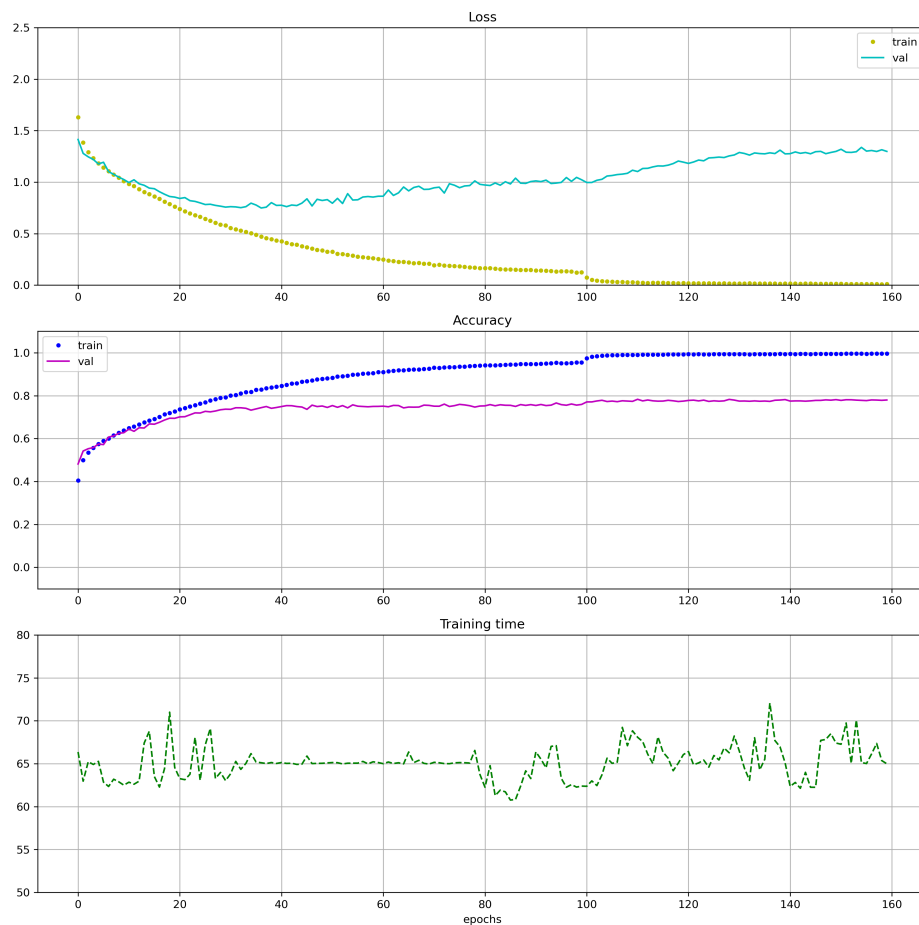
Listing 7: Funkcja wykonująca walidację po każdej epoce.

W poniższej pętli treningowej wykorzystałem obie powyższe funkcje, oraz zastosowałem metodę *Early stopping*.

```
1 vit_model = CifarViTClassifier().to('cuda')
2 vit_model.requires_grad_(True)
3
4 loss_function = nn.CrossEntropyLoss().cuda()
5
6 optimizer = AdamW(
7     vit_model.parameters(),
8     lr=5e-4
9 )
10
11 lr_scheduler = MultiStepLR(
12     optimizer=optimizer,
13     milestones=[100, 150],
14     gamma=0.1
15 )
16
17 train_data_loader = DataLoader(cifar_train_ds, batch_size=64, shuffle=True)
18 test_data_loader = DataLoader(cifar_test_ds, batch_size=64, shuffle=True)
19
20 train_loss_acc_time_list = []
21 val_loss_acc_list = []
22
23 max_test_acc = 0.
24 early_stopped_vit_model = CifarViTClassifier().to('cuda')
25 early_stopped_vit_model.load_state_dict(vit_model.state_dict())
26
27 n_epochs = 160
28
29 for epoch in range(n_epochs):
30     avg_train_loss, avg_train_acc, epoch_train_time = train_step(
31         vit_model,
32         train_data_loader,
33         loss_function,
34         optimizer
35     )
36
37     avg_val_loss, avg_val_acc = val_step(
38         vit_model,
39         test_data_loader,
40         loss_function
41     )
42
43     lr_scheduler.step()
44
45     train_loss_acc_time_list.append((avg_train_loss, avg_train_acc, epoch_train_time))
46     val_loss_acc_list.append((avg_val_loss, avg_val_acc))
47
48     if avg_val_acc > max_test_acc:
49         max_test_acc = avg_val_acc
50         early_stopped_vit_model.load_state_dict(vit_model.state_dict())
51
52     if epoch in (10, 25, 50, 100, 125):
53         torch.save(vit_model.state_dict(), Path.cwd() / f'vit-cifar10-{epoch}.pt')
```

Listing 8: Główna pętla treningowa

Historię treningu przedstawiłem na poniższym wykresie. Wyraźnie widać poprawę dokładności na zbiorze treningowym (a na walidacyjnym może i też) po setnej epoce (wtedy, kiedy `lr_scheduler` pomnożył *learning rate* przez 0.1).



Rysunek 3: Statystyki treningu

Przetestowałem model na kilku obrazach ze zbioru testowego i zazwyczaj, jeżeli ja byłem w stanie rozpoznać obraz, to sieć też sobie radziła dobrze (przykładowo: niektóre żaby są strasznie rozmazane i ciężko je rozpoznać :)).

## 6 Attention rollout

Aby uzyskać mapę atencji metodą *Attention rollout*, wykorzystałem funkcję, modyfikującą *forward hook* modułów **MultiHeadAttention**, tak aby zapisywane zostały macierze atencji. Następnie, przemnożyłem rekurencyjnie te macierze i uzyskałem *attention rollout matrix*.

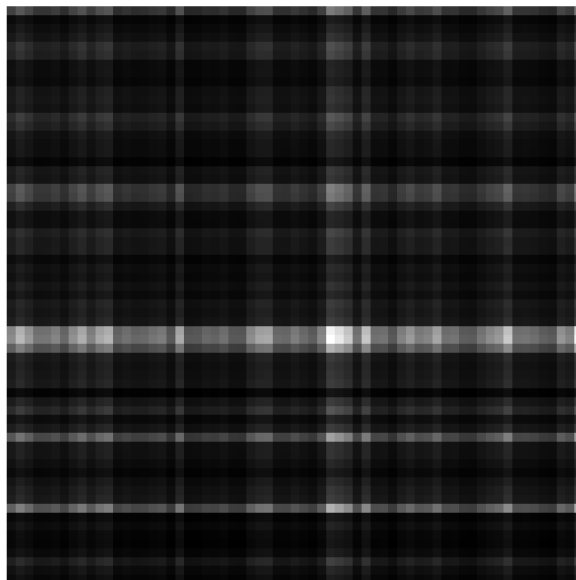
```
1 def get_mha_matrix_from_nth_transformer(  
2     vit_m: CifarViTClassifier,  
3     n: int,  
4     image: torch.Tensor  
5 ) -> torch.Tensor:  
6     """Get multi-head attention matrix from nth transformer using a forward hook."""  
7  
8     result_activation = None  
9  
10    def output_hook(model, _input, _output_tuple):  
11        nonlocal result_activation  
12        result_activation = _output_tuple[1].detach()  
13  
14    handle = vit_m.transformers_seq[n].mh_attention.register_forward_hook(output_hook)  
15  
16    _ = vit_m(image.unsqueeze(0))  
17  
18    handle.remove()  
19  
20    return result_activation.squeeze()
```

Listing 9: Funkcja zwracająca macierz atencji z n-tego transformera.

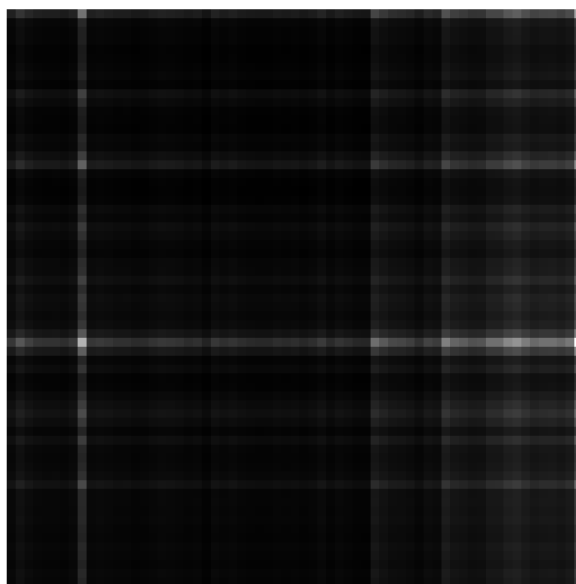
```
1 def get_patches_attention_matrix_rollout(  
2     vit_m: CifarViTClassifier,  
3     image: torch.Tensor  
4 ) -> torch.Tensor:  
5  
6     att_mat = get_mha_matrix_from_nth_transformer(vit_m, 0, image)  
7     for i in range(1, vit_m.n_transformers):  
8         ith_att_mat = get_mha_matrix_from_nth_transformer(vit_m, i, image)  
9         att_mat = att_mat.T @ ith_att_mat  
10  
11    return att_mat
```

Listing 10: Funkcja zwracająca attention rollout matrix (liczba patchy + 1 x liczba patchy + 1)





Rysunek 4: Macierz uwagi dla zdjęcia ptaka ze zbioru treningowego.



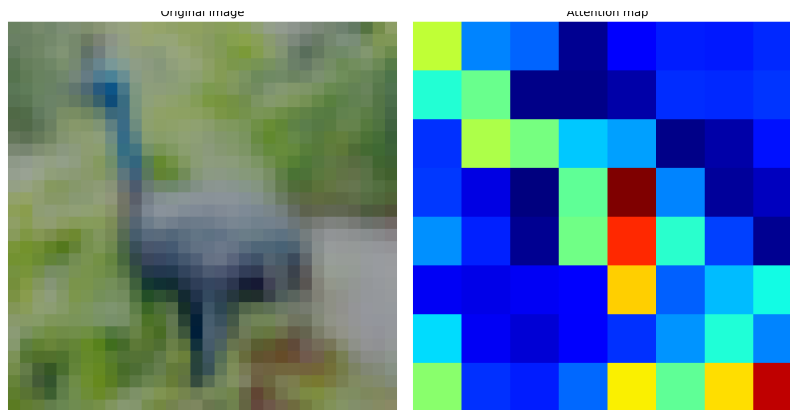
Rysunek 5: Macierz uwagi dla zdjęcia żaglówki ze zbioru testowego.

Mapy atencji uzyskałem z macierzy atencji korzystając z poniższego skrawku kodu.  
 Poniższe mapy przedstawiają (a przynajmniej powinny w teorii) atencję odnośnie *cls\_token*.

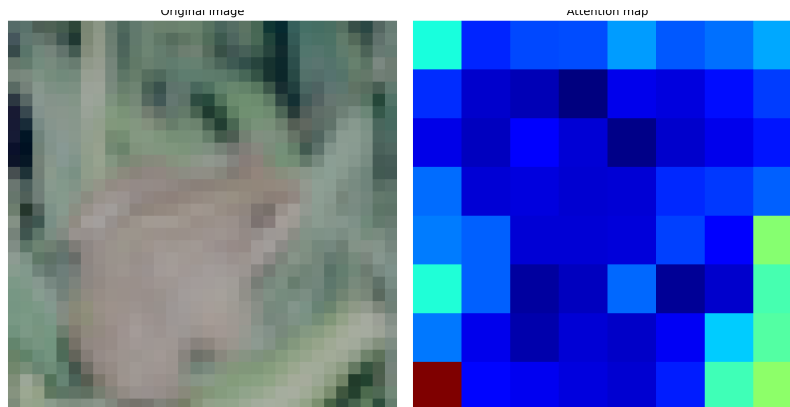
```

1 # plot attention with regard to cls token
2 enlarged_attention_map = att_matrix_rollout[0, 1:].reshape(
3     int(np.sqrt(n_patches)), int(np.sqrt(n_patches))
4 )
5 enlarged_attention_map = resize(
6     enlarged_attention_map.unsqueeze(0), size=(im_size, im_size),
7     interpolation=torchvision.transforms.InterpolationMode.NEAREST
8 ).squeeze()
9 plt.imshow(np.moveaxis(enlarged_attention_map.detach().numpy(), 0, -1), cmap='jet')
```

Listing 11: Skrawek prezentujący otrzymanie mapy atencji z attention matrix rollout i wyplotowanie jej.



Rysunek 6: Mapa atencji dla zdjęcia ptaka ze zbioru treningowego.



Rysunek 7: Mapa atencji dla zdjęcia żaby ze zbioru testowego.

Chętnie bym się dowiedział, czy zaimplementowałem metodę *attention rollout* poprawnie, bo mam podejrzenia, że gdzieś mógł wkraść się błąd. Druga możliwość jest taka, że sieć jest nadmiernie dopasowana do zbioru treningowego i należałoby poprawić generalizację modelu (np. stosując bardziej zaawansowane metody augmentacji).