

Computer Vision Assignment 2

Authors (Group 2):

- Aditya Shourya (i6353515)
- Ashkan Karimi Saberi (i6365010)

code at : https://github.com/adishourya/vit_fer2013

Structure of the Report

- The focus of the report would not be on attaining the best test score but would be on experimenting with different architectures.
- And looking through how the forward pass looks in each case.
- And we will try to reason rate of learning whenever applicable
- And finally we will pass in a pair of images as sampled from a video

Convolutional Neural Network

- we will first develop a simple convolutional neural network
- And reason the effectiveness of the network on a simple task as FER 2013 [<https://www.kaggle.com/datasets/msambare/fer2013>]

A simple hand made vision transformer (shallow)

- we will try to replicate the paper "AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE" -- dosovitskiy et.al [<https://arxiv.org/pdf/2010.11929v2>] and show through one full forward pass.

Pretrained Vision Transformer

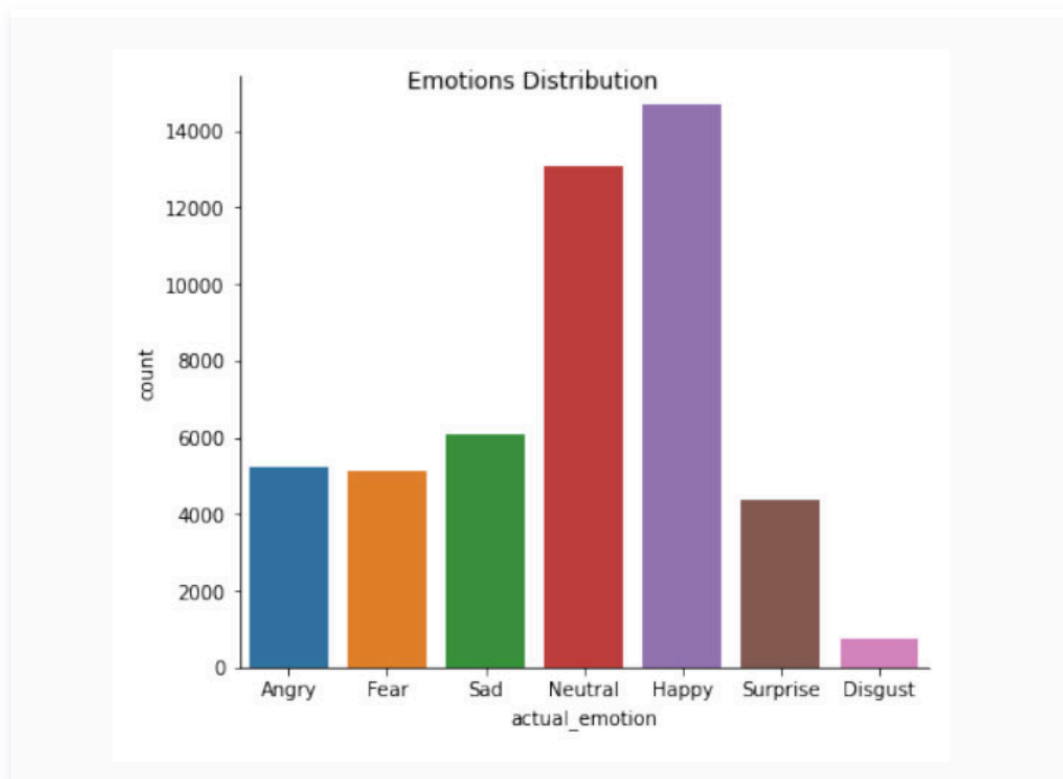
- we will use vit_b_16 [https://pytorch.org/vision/main/models/generated/torchvision.models.vit_b_16.html#vit-b-16] and unfreeze the last few layers and attention heads to perform transfer learning.

FER2013 Dataset

The dataset comprises **48x48** pixel grayscale images of faces. These faces have been automatically aligned to ensure that each face is roughly centered and occupies a similar amount of space in every image.

The objective is to classify each facial expression into one of seven emotional categories: 0 for Angry, 1 for Disgust, 2 for Fear, 3 for Happy, 4 for Sad, 5 for Surprise, and 6 for Neutral. The training set includes 28,709 examples, while the public test set contains 3,589 examples.

-- [<https://www.kaggle.com/datasets/msambare/fer2013>]



Class Distribution in FER2013, and the challenges associated with the dataset are :

The FER2013 dataset has several inherent issues that make it challenging for deep learning architectures to achieve optimal results. Key problems include **imbalanced data**, intra-class variation, and occlusion. Specifically, the dataset exhibits significant imbalance in the training data, with classes having vastly different numbers of samples. For example, the 'happy' emotion has over 13,000 samples, whereas 'disgust' has only about 600 samples, as shown in the figure above.

Intra-class variation refers to the differences within the same class. Reducing intra-class variation while increasing inter-class variation is crucial for effective classification. Variations, uncontrolled lighting conditions, and occlusions are common issues that face recognition systems encounter in real-world applications. These challenges often lead to a drop in accuracy compared to performance in controlled experimental settings. Occlusion occurs when an object blocks part of a person's face, such as a hand, hair, cap, or sunglasses. Although occlusion complicates face recognition, it can also provide useful information, as people often use their hands while communicating through gestures.

-- <https://www.oaepublish.com/articles/ir.2021.16>

CNN

- All the experiments and the code are from the notebook : https://github.com/adishourya/VIT_FER2013/blob/main/convolutional_augmentation.ipynb
- Basic Augmentations applied: horizontal and vertical flipping for our model to become invariant to geometrical transformations.
- We don't use any adjustment sharpness ; because we saw it made the lips and the expression overly smooth to discern emotion.

```
1 self.transforms = transforms.Compose([
2     transforms.ToTensor(),
3     transforms.RandomHorizontalFlip(p=0.5),
4     transforms.RandomVerticalFlip(p=0.5),
5     # transforms.RandomAdjustSharpness(sharpness_factor=2), #
    makes lips look worse and it performs slightly worse with it
6     transforms.Resize((32,32), antialias=True)]) # limit number
    of patches (4) calculation.. keep this a multiple of 16*16
7
```

- A simple neural network from the documentation in pytorch.

```

1  # define a small convolutional network
2  # see beautiful mnist in tinygrad .
3  import torch.nn as nn
4  import torch.nn.functional as F
5
6  # shape after operations n,n →(with padding p and stride s) (n + 2p - f +
  1 )/s + 1
7
8  class Net(nn.Module):
9      def __init__(self):
10         super().__init__()
11         self.conv1 = nn.Conv2d(1, 6, 5) # input channel 1 , 6 filter banks
  each of kernels size (5,5)
12         self.pool = nn.MaxPool2d(2, 2) # this is not a learnable operation
  just performs downsampling
13         self.conv2 = nn.Conv2d(6, 16, 5) # 16 kernels
14         self.fc1 = nn.Linear(16 * 5 * 5, 120)
15         self.fc2 = nn.Linear(120, 84)
16         self.fc3 = nn.Linear(84, 7) # we have 7 classes
17
18     def forward(self, x):
19         x = self.pool(F.relu(self.conv1(x)))
20         x = self.pool(F.relu(self.conv2(x)))
21         x = torch.flatten(x, 1) # flatten all dimensions except batch
22         x = F.relu(self.fc1(x))
23         x = F.relu(self.fc2(x))
24         logits = self.fc3(x)
25         return logits
26
27
28  net = Net()

```

- we choose commonly used 5,5 filters in our case. with padding of 0 and stride of 1.

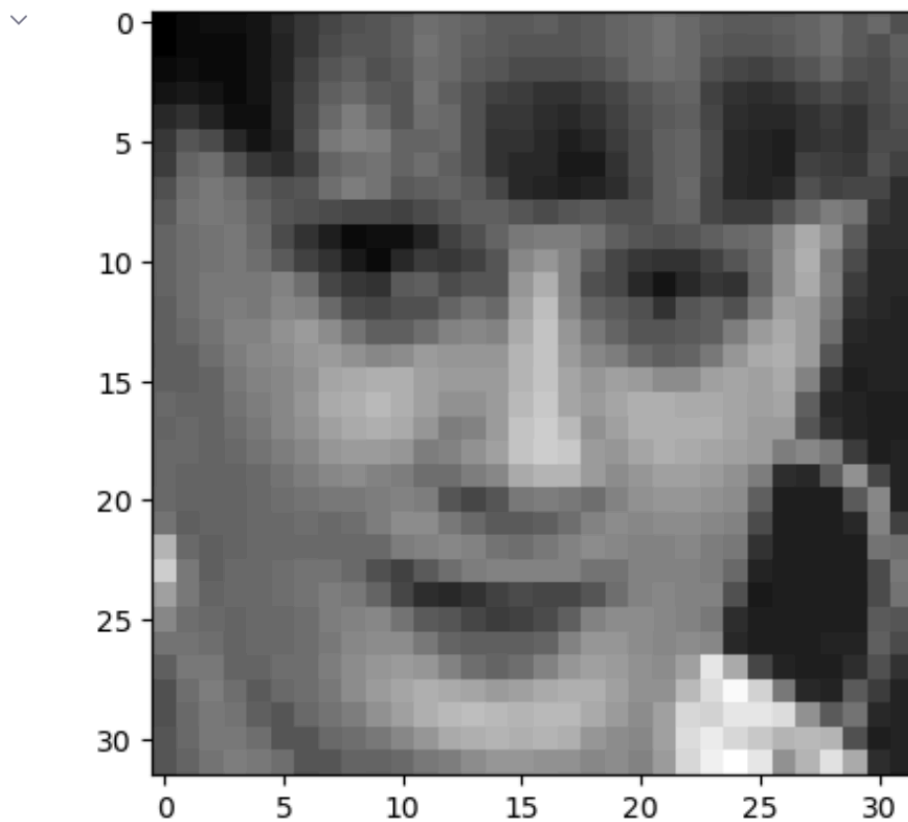
forward pass in our learnt convolutional layers

```
1  # lets visualize the first layer on a sample image
2  x = self.pool(F.relu(self.conv1(x)))
3
4  # 1,6 input channel = 1 (greyscale)
5  # output dimension = 6 (filterbanks)
6  # kernel size = (5,5) kernels
7  net.conv1 , net.conv1.weight.shape
8
9  # > out: (Conv2d(1, 6, kernel_size=(5, 5), stride=(1, 1)), torch.Size([6, 1,
    5, 5]))
```

```
In 52 1 import einops
      2 sample_img = img[0]
      3 plt.imshow(
      4     einops.rearrange(sample_img, "1 h w → h w"),
      5     cmap="grey")
```

Executed at 2024.05.24 22:29:07 in 247ms

Out 52 <matplotlib.image.AxesImage at 0x15796c460>



```
1 import torch.nn.functional as F
2 with torch.no_grad():
3     # first layer convolution
4     out1 =F.conv2d(sample_img,net.conv1.weight,
5                     bias=None, stride=1, padding=0)
6     print(out1.shape) # input channel = 1 filter banks = 6
7
8
9 plt.imshow(
10     einops.rearrange(out1,"out_c h w → h (out_c w)"),
```

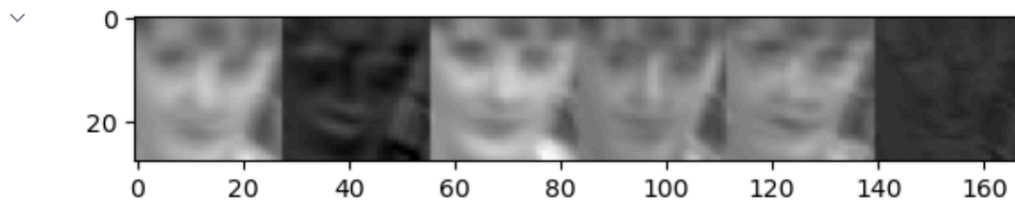
```

11     cmap="grey"
12 )
13 # visualizing the output of convolution from the first learnt layer.
14 # notice how some kernels are vastly different

```

```
torch.Size([6, 28, 28])
```

Out 63 <matplotlib.image.AxesImage at 0x157e77550>



- we will not show the pooling operation here (check notebook) as it is not a learnable parameter

```

1  with torch.no_grad():
2      # first layer convolution
3      out2 =F.conv2d(out1_p,net.conv2.weight,
4                      bias=None, stride=1, padding=0)
5      print(out2.shape) # input channel = 1 filter banks = 6
6
7  plt.figure(figsize=(15,8))
8  plt.imshow(
9      einops.rearrange(out2,"out_c h w → h (out_c w)"),
10     cmap="grey"
11 )

```

```
torch.Size([16, 23, 23])
```

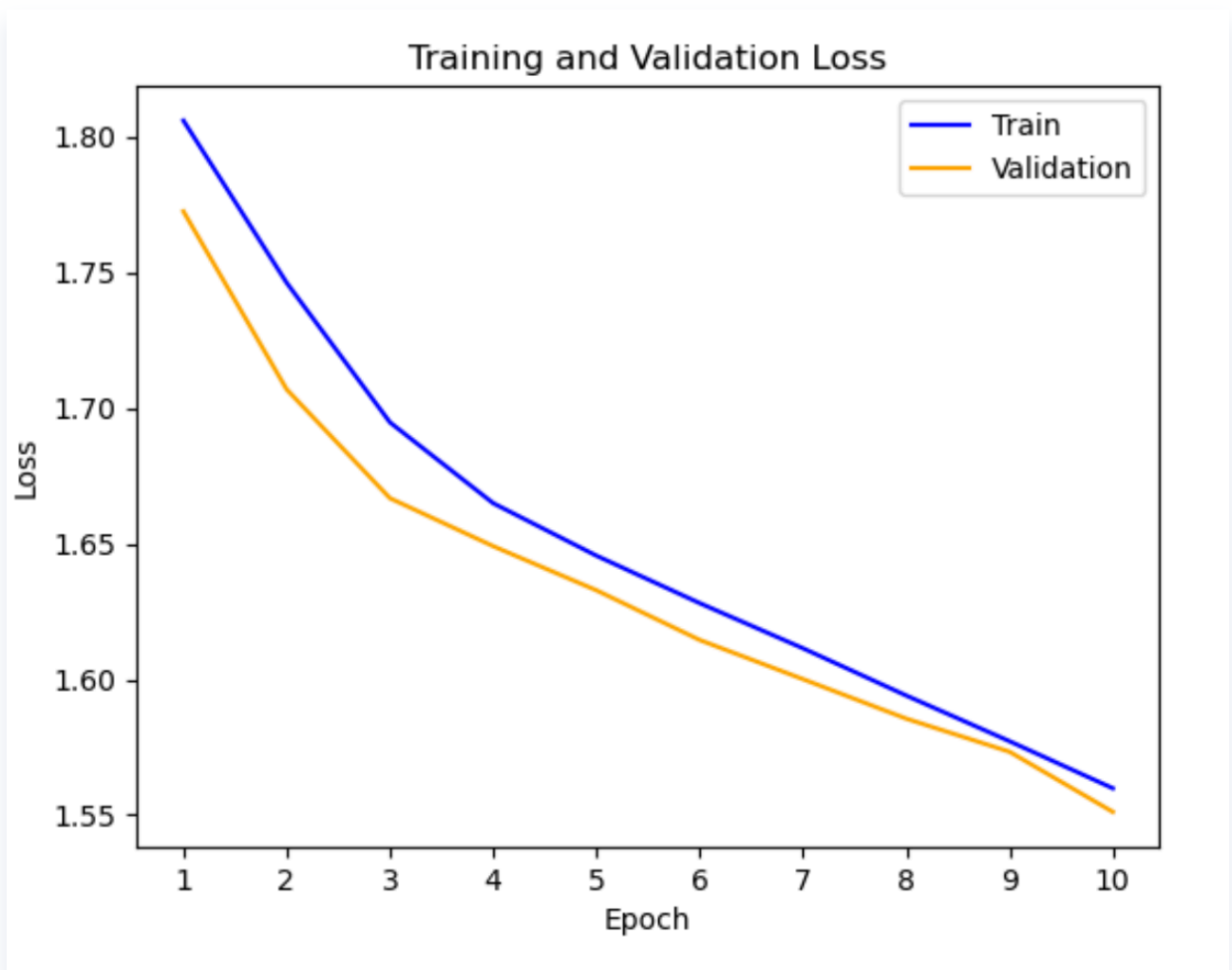
<matplotlib.image.AxesImage at 0x157fcbdc0>



- and then we maxpool it and then flatten it to pass to a feed forward network to arrive at the logits

for classification

- we only train for 10 epochs as attaining the best test accuracy is not our main goal of this report



Executed at 2024.05.24 21:24:47 in 25.129ms

100% | ██████████ | 113/113 [00:02<00:00, 53.20it/s]

Test Loss: 1.556882


```
1 Test Accuracy of Classes (Generalization scores of our Convolutional Model)
2
3 Angry : 10% (54/491)
4 Disgust : 0% (0/55)
5 Fear : 7% (41/528)
6 Happy : 80% (710/879)
7 Sad : 32% (191/594)
8 Surprise : 47% (197/416)
9 Neutral : 33% (208/626)
10
11 Test Accuracy of Dataset: 39% (1401/3589)
12
```

Shallow Transformer

- we will try to replicate a shallow vit from the paper dosovitskiy et.al [<https://arxiv.org/pdf/2010.11929v2>]
- AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

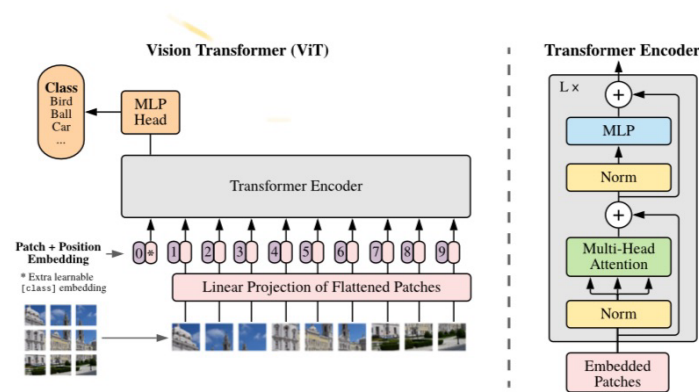


Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable “classification token” to the sequence. The illustration of the Transformer encoder was inspired by Vaswani et al. (2017).

3 METHOD

$$N = \frac{HW}{P^2} = \frac{32 \times 32}{16 \times 16} = 4$$

| | |
|-------|-------|
| P_1 | P_2 |
| P_3 | P_4 |

$$P \in \mathbb{R}^{16,16,C=1}$$

- All the experiments and the code are from : https://github.com/adishourya/ViT_FER2013/blob/main/fer_vit.ipynb
- we will first start with splitting the image as patches

generating patches

To do so, we split an image into patches and provide the sequence of linear embeddings of these patches as an input to a Transformer. Image patches are treated the same way as tokens

-- Introduction adosovitskiy et.al

Naive application of self-attention to images would require that each pixel attends to every other pixel.

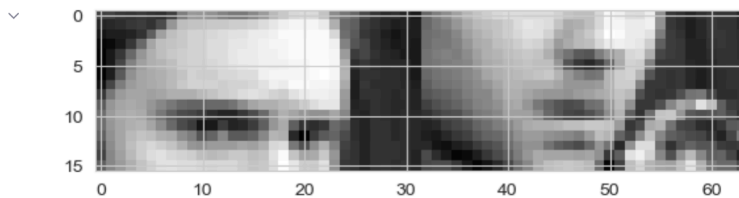
-- Related Work adosovitskiy et.al

```

1  ## from class FerDataset(utils.data.Dataset): (check Notebook fer_vit)
2  def get_patches(self, arr):
3      # best place to use einops
4      # squeeze out the channel dimension.. the loader will add it ..
5      arr = einops.rearrange(arr, "1 h w → h w")
6      patches = einops.rearrange(
7          arr, "(p1 h) (p2 w) → (p1 p2) h w", p1=2 , p2=2
8      ) # lay it on as the batch for the image
9      # we need to crop twice for a 48 * 48 image to get a 16 * 16 patch
10     return patches
11

```

Out 23 <matplotlib.image.AxesImage at 0x32c5ee530>



Embed Patches

- We first flatten all the patches and embed them with a lookup table

```

1  # initialize a lookup table
2  self.C = nn.Parameter(torch.randn(16*16, self.n_embed) * 1/torch.sqrt(256)
3  )
4
5  # flatten the patches
6  patches = einops.rearrange(X, "b p h w → b p (h w)")
7
8  # embed them : the block matrix multiplication looks like :
9  # B , p_num , 256 @ 256 , 32
10 emb = patches @ self.C # B , p_num , n_embed
11
12 # so each patch (a matrix) gets an n_embed dimensional representation (row
13 vector)
14

```

class token

- we dont actually use the class token in our implemented model. check screenshot of the issue :
` Is the extra class embedding important to predict the results, why not simply use feature maps to predict?` in the appendix.
- But the pretrained models implement them anyway. So we do it like :

```
1  # its like adding an extra patch . but this patch is learnable
2  learnable_patch = torch.randn(1,n_embed)
3  # add same learnable patch to all the images in the batch
4  learnable_patch = einops.repeat(learnable_patch,"1 n_embed → 3 1 n_embed")
5  print(learnable_patch.shape)
6  # prints torch.Size([3, 1, 128]) # batch_size , 1 token , 128 embed
   dimension
7
8  # you actually need all dims except one to be different to pack or concat
   it
9  # in the paper they add it at the first location. (dont really see how it
   matters if append it at the end instead)
10 xb ,ps = einops.pack([learnable_patch,xb],"b * n_embed") # we want to
    append it on top of patch
11 xb.shape
12 # prints torch.Size([3, 5, 128]) # batch , tokens , n_embed
13
14 # 4 patches from the images , 1 patch from class token
```

positional embedding

```

1  # positional embedding
2  # the paper says Position embeddings are added to the patch embeddings to
   retain positional information.
3  # We use standard learnable 1D position embeddings, since we have not
   observed significant performance gains from using more advanced 2D-aware
   position embeddings
4
5  # init a patch embedding
6  self.pe = nn.Parameter(torch.randn(1,4,self.n_embed)) # each pach and
   representation should get positional embedding
7
8  # add patch embedding after class token
9  emb = emb + self.pe # kind of acts like a bias towards each patches.

```

pass through transformer blocks

- We will first start with making self attention blocks. and then concatenate them to get multihead attention (vasawani et.al 2017) .
- Note multihead attention blocks returns back the same size of representation `n_embed` .
- unit test for our implementation of self head attention is in the appendix

```

1  # check appendix for unit test of this class
2  class SingleHead(nn.Module):
3      """
4      Implements a single head of attention (unmasked)
5      """
6
7      def __init__(self,n_embed=32,head_size=8):
8          super().__init__()
9          # single head
10         self.head_size = torch.tensor(head_size)
11         self.n_embed = torch.tensor(n_embed)
12         self.Q = nn.Parameter( torch.randn(self.n_embed,head_size) *
(1/torch.tensor(2.82)))
13         self.K = nn.Parameter( torch.randn(self.n_embed,head_size) *
(1/torch.tensor(2.82)))
14         self.V = nn.Parameter( torch.randn(self.n_embed,head_size) *
(1/torch.tensor(2.82)))
15

```

```

16     def forward(self,x):
17         query = x @ self.Q
18         key = x @ self.K
19         value= x @ self.V
20
21         # hand implementation
22         # scale  $\Rightarrow$  sqrt head size
23         scale = 1 / torch.sqrt(self.head_size)
24
25         # we will not use any masking here as its an image
26         # and no dropout consideration in this implementation
27         comm = query @ key.transpose(-2,-1)
28         comm = comm* scale
29         soft_comm = torch.softmax(comm, dim=2)
30         att = soft_comm @ value
31
32         return att

```

```

1  class Multihead(nn.Module):
2      def __init__(self,n_embed,n_heads):
3          super().__init__()
4
5          self.n_embed = n_embed
6          self.n_heads = n_heads
7          self.head_size = self.n_embed // self.n_heads
8
9          self.multiheads = nn.ModuleList(
10              [SingleHead(self.n_embed,self.head_size)
11               for _ in range(self.n_heads)]
12          )
13
14      def forward(self,x):
15          return torch.cat([head(x) for head in self.multiheads],dim=2)

```

- And then a transformer block with it

```

1  # only multihead  $\rightarrow$  skip connection  $\rightarrow$  layernorm
2  # Batch norm : couples examples in and normalizes it .. (also has a
    regularization effect) but we need to keep a running mean to track new mean
    and sigma

```

```

3 # layernorm : normalizes the features of each example (does not couple
  examples across the batch) more popular in transformers
4
5 class TranformerBlock(nn.Module):
6     def __init__(self, n_embed, n_head):
7         super().__init__()
8         self.multi_head = Multihead(n_embed, n_head)
9         # i am not going to implement my own layer norm it wont be
  efficient and will be janky at best
10        self.norm = nn.LayerNorm(n_embed) # we want to normalize ffeatures
  (each patch gets normalized)
11
12    def forward(self, x):
13        # pass through multihead
14        attention = self.multi_head(x)
15        # skip connection and non linearity
16        attention = torch.relu( x + attention)
17        # layer norm
18        attention = self.norm(attention)
19        return attention # B , n_patch , n_embed
20        ...
21

```

- So our architecure looks like this :

```

1 # most of the comments are pasted verbatim from the paper
2 class SmallVIT(nn.Module):
3
4     def __init__(self):
5         super().__init__()
6         # patches
7         # embedding
8         self.vocab_size = torch.tensor(256) # 0 to 255 pixels
9         # each patch will only get one n_embed representation
10        self.n_embed = 32 # we will project each patch 16*16 to a 32
  dimensional representation
11        # so the lookup table would be of the shape
12        # unlike in nlp where we embed token to a vector like below we
  would project matrix of patch size to a vector
13        # self.C = nn.Embedding(self.vocab_size, self.n_embed)
14

```

```

15         self.C = nn.Parameter(torch.randn(self.vocab_size, self.n_embed) *
16                                   1/torch.sqrt(self.vocab_size) )
17
18         # positional embedding
19         # the paper says Position embeddings are added to the patch
20         embeddings to retain positional information. We use standard learnable 1D
21         position embeddings, since we have not observed significant performance
22         gains from using more advanced 2D-aware position embeddings
23         self.pe = nn.Parameter(torch.randn(1,4,self.n_embed)) # each patch
24         and representation should get positional embedding
25
26         # we use the standard approach of adding an extra learnable
27         "classification token" to the sequence
28         self.classification_token = nn.Parameter(torch.randn(1, 1,
29 self.n_embed))
30
31         # we will keep the step above optional .. i dont understand why we
32         should use it yet.
33
34         # transformer block
35         self.n_heads = 4 # we will use 4 heads for now
36         self.transformer_block = TransformerBlock(self.n_embed,self.n_heads)
37
38         # MLP Head for final logit calculation
39         # n_patch * n_embed → fer["emotion"].nunique() : 7
40         self.mlp_head = nn.Parameter(torch.randn(4*32, 7) *
41 torch.sqrt(torch.tensor(4*32)))
42
43
44
45
46
47
48         ...
49
50     def forward(self,X):
51         batch_size = X.shape[0]
52         # flatten the patches
53         patches = einops.rearrange(X,"b p h w → b p (h w)")
54         # B , p_num , 256 @ 256 , 32
55         emb = patches @ self.C # B , p_num , n_embed
56         emb = emb + self.pe # kind of acts like a bias towards each
57         patches.
58
59         # 2 transformers
60         tf = self.transformer_block(emb)

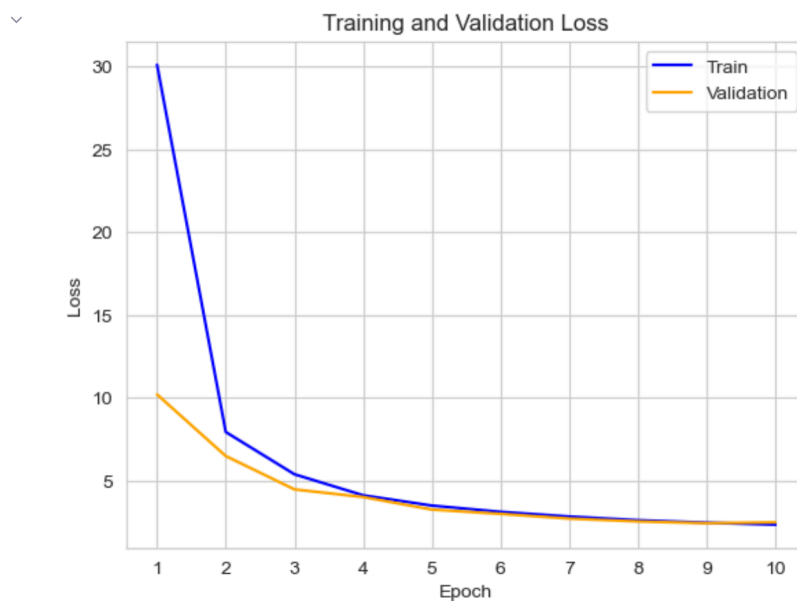
```



```

48         tf = self.transformer_block(tf)
49
50         # flatten it : across patches
51         tf = tf.view(batch_size,-1)
52
53         # logits
54         logits = tf @ self.mlp_head
55
56
57
58
59         # broadcasting steps in the above command
60         # B,p_n , p*p , n_embed
61         #1,4,32
62
63
64         return logits

```



100% |██████████| 113/113 [00:02<00:00, 54.60it/s]

Test Loss: 2.496379

```

1  Test Accuracy of Classes
2
3  Angry : 12%      (61/491)
4  Disgust : 0%     (0/55)
5  Fear   : 31%     (168/528)
6  Happy  : 62%     (551/879)
7  Sad    : 11%     (68/594)
8  Surprise : 13%   (56/416)
9  Neutral : 10%    (64/626)
10
11 Test Accuracy of Dataset: 26% (968/3589)

```

- Note how the accuracy of a shallow transformer (significantly higher number of parameters compared to CNN) produces worse result. which is to be expected.
- This looks like its improving but even after 10 epochs the model is still confidently wrong.

```

1  # if we were to assume the model is predicting randomly from a uniform
    distribution
2  # i.e logits were roughly equal for all classes, then the loss would have
    been :
3  -1 * np.log(1/7) # which is 1.946. which is still lower than what we have
    after 10 epochs!

```

- note how other models (both cnn and pre-trained vit) already has a loss score of < 1.946 after first epoch.
- This implies that the model was not initialized well , and since the model plateaus around the half way . it also indicates that the model does not have enough representational capacity.

Pretrained Vision Transformer (perform transfer learning)

- The code and experiments for this section can be found at : https://github.com/adishourya/VIT_ER2013/blob/main/vit_pretrained.ipynb
- we will be transfer learning from the model : vit_b_16 [https://pytorch.org/vision/main/models/generated/torchvision.models.vit_b_16.html#vit-b-16] and unfreeze the last few layers to perform the training.

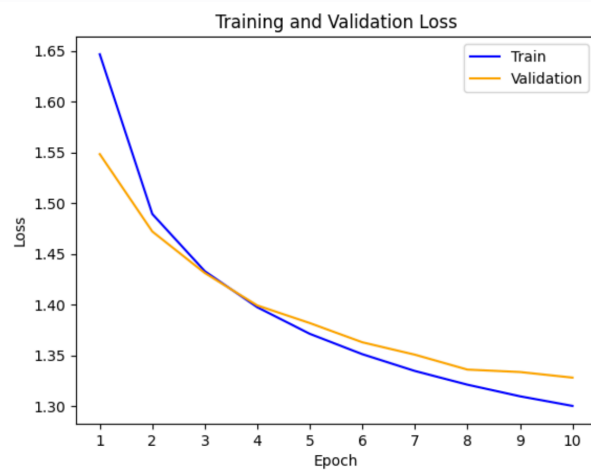
```
1 vision_transformer =  
  models.vit_b_16(weights=models.ViT_B_16_Weights.DEFAULT)
```

```
1 vision_transformer.heads  
2  
3 #Sequential(  
4 #   (head): Linear(in_features=768, out_features=1000, bias=True)  
5 #)
```

Fine Tune

```
1 # fine-tune with dataset  
2  
3 # change the number of output classes  
4 vision_transformer.heads = nn.Linear(in_features=256*3, out_features=7,  
  bias=True)  
5  
6 # freeze the parameters except the last linear layer  
7 #  
8 # freeze weights  
9 for p in vision_transformer.parameters():  
10     p.requires_grad = False  
11  
12 # unfreeze weights of classification head to train  
13 for p in vision_transformer.heads.parameters():  
14     p.requires_grad = True
```

- and then train



```
52
53 # output test loss statistics
54 print('Test Loss: {:.6f}'.format(test_loss))
55
```

100%|██████████| 113/113 [36:49<00:00, 19.55s/it]

Test Loss: 1.314458

```

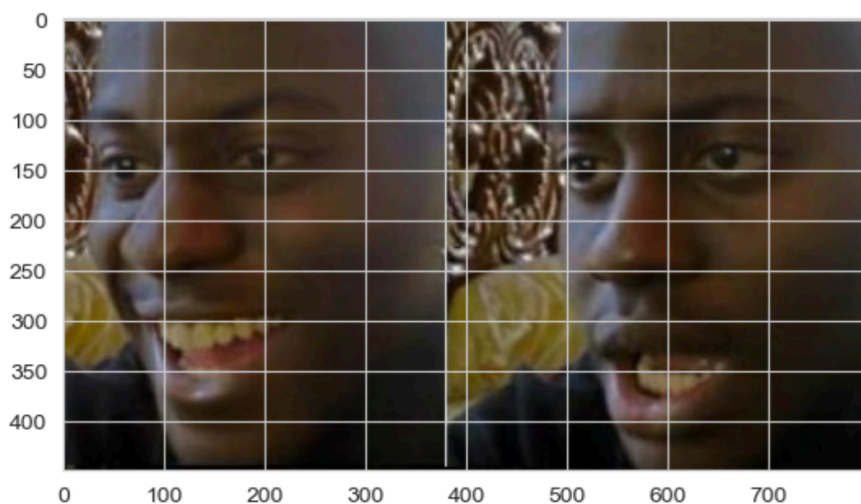
1  ## Results for each class
2  Disgust : 1%    (1/55)
3  Fear   : 29%   (155/528)
4  Happy  : 78%   (688/879)
5  Sad    : 35%   (211/594)
6  Surprise : 62% (260/416)
7  Neutral : 54%  (344/626)
8
9  Test Accuracy of Dataset:    50%  (1799/3589)

```

- Note how its considerably better than shallow transformer that we made
- It also performs better than our small convolutional network (Note that the CNN has extremely lower number of parameters compared to vit_b16)

Testing on images sampled from a real video

- we test the images on our CNN model.



- transform the image as expected by the model

```

1  # convert to grey scale
2  our_img= einops.reduce(our_img,"h w c → h w","min")
3
4  # crop image into batch size of 2
5  # 2 batches
6  ims = einops.rearrange(our_img," h (p w) → p h w",p=2)

```

```

7
8 # resize to the same shape as the training data and unsqueeze the channel
  dimension
9 ims= transforms.Resize((32,32), antialias=True)(ims)
10 ims = einops.rearrange(ims,"b h w → b 1 h w")
11
12 # forward pass through our network.
13 with torch.no_grad():
14     logits = net(ims)# logits
15     print("logits",logits)
16     probs = torch.softmax(logits,dim=1)
17     print("probs",probs)
18     print(torch.argmax(probs,dim=1))

```

```

1 logits tensor([[ 8.8504, -107.3212,  49.5044, -21.4002, -39.1525,
2                -52.2730],
3                [ 25.8041, -122.4651,  48.1215, -10.1948, -35.4007,  67.4436,
4                -51.2039]])
5 probs tensor([[4.4226e-36, 0.0000e+00, 2.0023e-18, 0.0000e+00, 0.0000e+00,
6                0.0000e+00],
7                [8.2455e-19, 0.0000e+00, 4.0601e-09, 1.9146e-34, 2.8026e-45,
8                1.0000e+00,
9                0.0000e+00]])
10 tensor([5, 5])
11 classes[5] → "surprise" # they were both classified as surprise.
12 # which is only partially correct . correct for the second image

```

Remarks

- Number of parameters:
 - CNN uses significantly lower number of parameters because of the inherent nature of parameter sharing in convolutional neural network. And performs comparably to a vision transformer that's not been trained for a lot of epochs.
 - when given enough scale , attention blocks could start learning context like the structural

pattern as the filter banks do in convolutional networks.

- Discuss difficulty of task at hand
 - The user should use a simple convolutional network when given a task with less sample size.

Transformers lack some of the inductive biases inherent to CNNs, such as translation equivariance and locality, and therefore do not generalize well when trained on insufficient amounts of data. -- adosovitskiy et.al

However, the picture changes if the models are trained on larger datasets (14M-300M images). We find that large scale training trumps inductive bias. Our Vision Transformer (ViT) attains excellent results when pre-trained at sufficient scale and transferred to tasks with fewer datapoints. When pre-trained on the public ImageNet-21k dataset or the in-house JFT-300M dataset, ViT approaches or beats state of the art on multiple image recognition benchmarks.

-- - adosovitskiy et.al

Appendix

unit test for self attention with pytorch's implementation

```
1 def single_head(query, key, value):
2     head_size = torch.tensor(query.shape[-1])
3     # hand implementation
4     # scale  $\Rightarrow$  sqrt head size
5     scale = 1 / torch.sqrt(head_size)
6
7     # we will not use any masking here as its an image
8     # and no dropout consideration in this implementation
9     comm = query @ key.transpose(-2, -1)
10    comm = comm * scale
11    soft_comm = torch.softmax(comm, dim=-2)
12    att = soft_comm @ value
```

```

13     print(att.shape)
14     return att
15
16
17     g=torch.Generator().manual_seed(123)
18     query, key, value = torch.randn(2, 3, 8 , generator = g), torch.randn(2, 3,
19     8, generator = g), torch.randn(2, 3, 8 , generator = g)
20
21     # our implementation
22     sh = single_head(query,key,value)
23
24     # pytorch implementation
25     py_sa = nn.functional.scaled_dot_product_attention(query, key, value)
26
27     # > torch.allclose(py_sa , sh) prints True

```

importance of class tokens ?

Track issue at : https://github.com/google-research/vision_transformer/issues/61

Is the extra class embedding important to predict the results, why not simply use feature maps to predict? #61 New issue

Closed QiushiYang opened this issue on Jan 26, 2021 · 3 comments

QiushiYang commented on Jan 26, 2021

Different from the common ways to use feature maps to obtain classification prediction (with fc or GAP layers), ViT employs an extra class embedding to do this without using feature maps explicitly. Wonder the meanings of this unusual design?

BTW, I used official pre-training params to fine-tune ViT on a small dataset, found that the validation accuracy is a little better after I replaced the feature maps with learnable class embedding to predict. So is the class embedding (maybe like a kind of query within encoder) important to learn and to predict?

lucasb-eyer commented on Mar 18, 2021 Collaborator

Great question. It is not really important. However, we wanted the model to be "exactly Transformer, but on image patches", so we kept this design from Transformer, where a token is always used.

Assignees: No one assigned

Labels: None yet

Projects: None yet

Milestone: No milestone

Development: No branches or pull requests

Notifications: Customize

- convolutional network without augmentation (which performs better) is at : https://github.com/adishourya/ViT_FER2013/blob/main/convolutional.ipynb

References

- Dataset used for experiment : FER 2013 [<https://www.kaggle.com/datasets/msambare/fer2013>]

- Tried to reproduce the paper : **AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE** -- adosovitskiy et.al [<https://arxiv.org/pdf/2010.11929v2>]
- Performed Transfer learning with : vit_b_16 [https://pytorch.org/vision/main/models/generated/torchvision.models.vit_b_16.html#vit-b-16]