

On a higher level they replace the need for domain or single label on image label annotation in Image-to-Image translation by training a guiding network that is able to do a self clustering of the image domain and therefore that guides the image-to-image translation, instead of the previously needed labels.

Until this paper these G-G translation requires labels because you need to know how to build these domains at the top. to get the different style vectors out & we will actually need label annotation for each image ie for each single image you need to know which one of the source corresponds to which one of the target.

~~Stages~~

Corresponding Images \rightarrow Sketch of a shoe

meaning the model does not require a labelled pair of images to train the network.

Stages that 2-2 translation went through

1) Corresponding Images One-to-one

Here is an example of the sketch of the shoe and here is the ~~corresponding~~ shot,

and similarly with the other shoes..

And from that you could learn a model that translates from one domain to the other because you have corresponding image-level annotations. (~~which~~ basically which element of domain A corresponds to which element of domain B.)

2) Set-level Annotations → That is basically the supervised labels for domain

→ domains. These things are instances of class A and shoe of class B. Yet there's no correspondence b/w the 2 classes. so 1-1 translation is not possible b/w domains when we just have domain level labels.

3. This paper introduced the following where you do not have domains any more but you just have the data set. Now the hypothesis is that there are still going to be domains in the dataset. They can be overlapping or not but there are still going to be domains, we just don't know what they are. Here there can be many ways to differentiate e.g.民族, hair color etc. So in essence we are going to assume that there are still some kind of a domain structure you just don't know what it is, but if we knew what it was then we could simply apply methods from set-level and the work will be done.

Now, what this paper shows is that if we apply something like a self clustering approach on x and other from basic I-I translation then this method is not as good as of quality the quality will not be as good as if we do both the things jointly.

So, what this paper does is that it jointly learns to cluster (i.e. self label) the I-maps and to do I-I translation. And by doing these jointly they help each other perform better.

Now how ~~we~~ do they do this?

There are 3 different parts to their model.

1. Encoder → which is the guiding network
2. Generator } So, the generator and discriminator
3. Discriminator } they are the standard GAN generators and discriminators i.e.

General Adversarial Network, but they do have some changes.

First We will see the discriminator.

It gets an image which can either be a generated or a real image and it needs to decide that whether the image is real or fake.

~~Here the input is a vector or image in this~~

Here we can see that input domain is a vector.

~~or for this. However, there are multiple heads~~

This whole thing is built on the pseudo-clustering approach where the pseudo labels which are coming from the left side here.

In essence we can assume that there are multiple classes in the dataset and the heads over here corresponds to each class. From somewhere outside, we will get the information of

CD O

E/F

the classification that what it's supposed to be, and then there is a head at the top of classifier which discriminates the data from our classification.

C

Hence, this discriminator is sort of a conditional discriminator - conditioned on a label. It

te it is just a labelled conditioned discriminator that discriminates your real and fake.

Now how we train the discriminator - that we would give it an image and we would let the encoder label the image (the process for which we will see after this) and for that particular label we will let the discriminator label ~~this~~ ~~as~~ ~~whether~~ clarify the image as real or fake.

Moons coming to the generator

The Encoder what it does is that it will take an image and it will output 2 things.

1) Label or style code

1

It's supposed

to be a member

~~b/v~~ 0 - k-1

~~is supposed to be~~

a class label → we just go and we sleep

These two things go on a different pathway
The label is directly going to the discriminator
~~the generator~~ does not see the label.

and the style code goes to the generator.
Now the generator takes a source image
and it takes the style code. Now the
style code is encapsulating the style of
the reference image.

The way we gonna train is that the style is
going to describe ^{the style} somehow all the images
that are from this label.

So the generator takes both the things and
produce a generated image. The discriminator
is then tasked with differentiating that whether
the image is real or fake for the given label.
over now.

And all of this is trained jointly. If we
jointly train the encoders to produce its
class labels and the styles. ^{we} train the generator
to take in the style & the source image and output an
image to feed the discriminator and the discriminator
at the same time is trained to differentiate b/w real and
fake images based on the label that the encoder gives.

At the end of all

→ we first need a way to train the encoder to come up with suitable class labels even though it does not get any back propagation signal into that part of its network.

~~8:12~~
So, now that's when we come to our loss functions. The way we are going to do this is that we are going to take the following approach: we are going to take an image x and we are going to take a randomly augmented version x^* . These are basically the ideas from self-supervision which is the work of some other paper.

We are going to take an image and we are going to augment it by random cropping, change its luminance etc. so we have two versions of the same image.

And we want to maximise the mutual information in not b/w the images themselves but p_v is going to be the output of the encoder
here

P.T.O
↔

If we say that

X goes into the encoder and the encoder outputs style and label class

+ noise

P is going to be the ~~(p)~~ class distribution of a histogram, from which we are going to sample the label Y

P is the ~~distribution~~ over the output classes

So, since we don't have a label we can't train the distribution. So what we will have to say is that we want to maximize the mutual information

$I(P)$

the Output Distribution of the image ($I_{P|I}$) and the output distribution of the augmented image ($I_{P|I'}$)

That entails the following two parts.

Entropy of P and the conditional entropy of P given P augmented

We want to maximize $\text{ent}(P)$ to minimize of MI.

In simple language what it means is that



we want different x 's let say $x_1, x_2 \dots, x_n$ and we want them to have different distributions in labels.

So, if the entropy is really high of the distribution p that means that different images get assigned to different classes. and if that is low that will mean that all the images get assigned to the single class.

So that's why the first thing that we want to maximise the entropy to maximize or the mutual information. and second we need to minimize the conditional entropy of p given p augmented.

Now, what it means is that if we know the augmented version of the image, its class labelling should be same as the unaugmented version; which means that if we take one of the x to x_i^+ , then that should keep the same class labelling.

e.g. a airplane in front of a blue sky is still an airplane in front of a bit blue sky. So that is why they should have the same output.

So by this Mutual Information we sort of ~~get~~ clustering of the output space.

By maximizing the MI, the network is encouraged to distribute all the samples as evenly as possible.

Page No. _____
Date: / / 201

That's how we train the pseudo clustering approach.

Now, we have a label to each image.

So, now how do we train the other parts.

The other parts are the additional losses like

→ π_i denotes the K -dimensional marginal probability vector

and π_{ij} denotes the joint probability.

Coming to the style Part.

Encoder outputs the labelling we got that covered and it outputs the style part. Now as we saw earlier the style part goes to the generator and in this part we also do get a generator back propagation to the style part; which implies that our encoder is trained to help the generator with its task of fooling the discriminator.

But first we will look at the loss function of the style loss of the encoder.

The style code is the contrastive loss $\frac{\text{numerator}}{\text{denominator}}$.

For this let us assume that ~~this~~ if we have our dataset and we take images out and train our dataset.

We take some n batches of images

so, we sort of have a queue. And we keep on pushing the new image into the queue and ~~throwing the old image out~~. So, for each image there are some other images. It is always going to be a queue of some other images.

The denominator over here is the others to the other images.

$s \rightarrow$ style code of the image you are considering right now

$s^+ \rightarrow$ style code of the augmented image.

~~s^-~~ Again considering X_1, X_2, X_3, \dots , we have our X , we put it through the encoder and we get our style code, and we augment X_1 to be X_1^+ and when we put that through the encoder it gives us s_1^+ and we put all the other images to the encoder and they are the s^- .

Now what we require is that the style code of our image is closer to the style code of its augmented version ~~than the previous~~ than it is to any of the other images style code.

So it's a combination of two things where you pull together things that are close to each other and you push apart things that we think should be far away from each other.

Here are some other loss functions used in the paper

- 1) The model uses the adversarial loss. This loss is not that different from the loss function used in typical GAN MODELS.
- The output D denotes the logic from the domain specific discriminator and S -TELDA denotes the target style code vector from the guidance network.

The generator G - learns to translate input image X to the target domain Y - TELDA while reflecting this style code.

In order to prevent degenerate situations where the generator ignores the given style code vector the model also use the style contrastive loss. The loss function guides the generated image to have a style similar to the reference image and dissimilar to random negative samples

Lastly to ensure that the generator can reconstruct its source image x when given with its original style s .

The model uses a reconstruction loss. The function not only ensures generator to preserve domain invariant characteristics like pose of its input image but also helps to learn the style representation of the guiding network by extracting the original style of the source image.

The final loss function is shown. The lambdas are the hyperparameters.

Coming to the overall loss function.

- 1) Discriminator \rightarrow The discriminator is just the adversarial loss
- 2) Generator \rightarrow has the Adversarial loss + The ^{style} content consistency + its own image content consistency.
- 3) Encoder \rightarrow All the generator losses + Mutual Info loss + Style contrastive loss

Comparisons

What if we separately train the guide network
and then beamlet. →

Result → no + SNE → non lens.

FIN → Quality matrix for GAN's lower is better.