
LEARNING PRIORS FOR ADVERSARIAL AUTOENCODERS

A PREPRINT

Belozerova Polina

Skoltech

bel.pol.4@gmail.com

Safin Alexander

Skoltech

safinsam@yandex.ru

Pavlovskaiia Natalia

Skoltech

ya-ne-bo@yandex.ru

October 26, 2018

ABSTRACT

Most deep latent factor models choose simple priors for simplicity, tractability or not knowing what prior to use. Recent studies show that the choice of the prior may have a profound effect on the expressiveness of the model, especially when its generative network has limited capacity. In this paper, we propose to learn a proper prior from data for adversarial autoencoders (AAEs). We introduce the notion of code generators to transform manually selected simple priors into ones that can better characterize the data distribution. Experimental results show that the proposed model can generate better image quality and learn better disentangled representations than AAEs in both supervised and unsupervised settings. Lastly, we present its ability to do cross-domain translation in a text-to-image synthesis task.

We took here several attempts to create a generative model. The main source of inspiration is [1]. The results are presented for MNIST and CIFAR-10 datasets. We haven't achieve the same quality as in the original paper, but learnt a lot.

1 Introduction

1.1 Goal

Our goal is to reproduce the "Learning priors for adversarial autoencoders" [1]

2 Algorithms

2.1 Adversarial Autoencoder

Our baseline model is "Adversarial Autoencoders" [2]

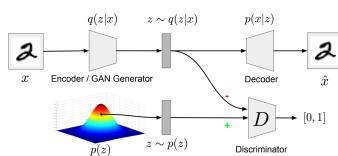


Figure 1: Adversarial Autoencoder scheme

2.2 Original paper algorithm

Algorithm 1 Training algorithm for our method.

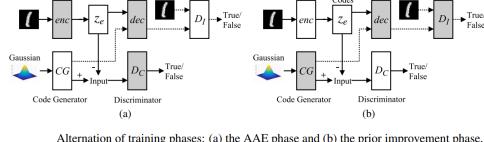
```

 $\theta_{enc}, \theta_{dec}, \theta_{CG}, \theta_{D_I}, \theta_{D_G}, \theta_Q \leftarrow$  Initialize network parameters
Repeat (for each epochs,  $E_i$ )
  Repeat (for each mini-batch  $x_j$ )
    // AAE phase
     $z \sim p(z)$ 
    If conditional variables  $s$  exist then
       $z_c \leftarrow CG(z, s)$ 
    Else
       $z_c \leftarrow CG(z)$ 
    End If
     $L_{GAN}^C \leftarrow -\log(D_C(z_c)) + \log(1 - D_C(enc(x)))$ 
     $x_{rec} \leftarrow dec(enc(x))$ 
     $L_{rec} \leftarrow \frac{1}{N} \| \mathcal{F}(x) - \mathcal{F}(x_{rec}) \|_2$ 
    // Update network parameters for AAE phase
     $\theta_{D_I} \leftarrow \theta_{D_I} - \nabla_{\theta_{D_I}} (L_{GAN}^C)$ 
     $\theta_{enc} \leftarrow \theta_{enc} - \nabla_{\theta_{enc}} (-L_{GAN}^C + L_{rec})$ 
     $\theta_{dec} \leftarrow \theta_{dec} - \nabla_{\theta_{dec}} (\lambda * L_{rec})$ 
    // Prior improvement phase
     $z \sim p(z)$ 
    If conditional variables  $s$  exist then
       $z_c \leftarrow CG(z, s)$ 
    Else
       $z_c \leftarrow CG(z)$ 
    End If
     $x_{noise} \leftarrow dec(z_c)$ 
     $x_{rec} \leftarrow dec(enc(x))$ 
     $L_{GAN}^C \leftarrow -\log(D_I(x_c)) + \log(1 - D_I(x_{noise})) + \log(1 - D_I(x_{rec}))$ 
    // Update network parameters for prior improvement phase
     $\theta_{D_I} \leftarrow \theta_{D_I} - \nabla_{\theta_{D_I}} (L_{GAN}^C)$ 
    If conditional variables  $s$  exist then
       $\theta_{dec} \leftarrow \theta_{dec} - \nabla_{\theta_{dec}} (-L_{GAN}^C + I(s; dec(z_c)))$ 
       $\theta_Q \leftarrow \theta_Q - \nabla_{\theta_Q} I(s; dec(z_c))$ 
    Else
       $\theta_{dec} \leftarrow \theta_{dec} - \nabla_{\theta_{dec}} (-L_{GAN}^C)$ 
    End If
    Until all mini-batches are seen
Until terminate

```

(a) Original paper algorithm pseudo code

Figure 2: Original paper algorithm



Alternation of training phases: (a) the AAE phase and (b) the prior improvement phase.

(b) Original paper algorithm scheme

2.3 Additional algorithm

Is inspired by "Autoencoding beyond pixels using a learned similarity metric"[3]

Algorithm 1 Training the VAE/GAN model

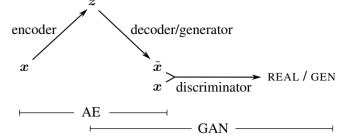
```

 $\theta_{Enc}, \theta_{Dec}, \theta_{Dis} \leftarrow$  initialize network parameters
repeat
   $X \leftarrow$  random mini-batch from dataset
   $Z \leftarrow Enc(X)$ 
   $L_{prior} \leftarrow D_{KL}(q(Z|X) || p(Z))$ 
   $\hat{X} \leftarrow Dec(Z)$ 
   $L_{Dis^{Dis}} \leftarrow -\mathbb{E}_{\hat{X}}[\text{Dis}(\hat{X}|Z)]$ 
   $Z_p \leftarrow$  samples from prior  $\mathcal{N}(0, I)$ 
   $X_p \leftarrow Dec(Z_p)$ 
   $L_{GAN} \leftarrow \log(\text{Dis}(X)) + \log(1 - \text{Dis}(\hat{X})) + \log(1 - \text{Dis}(X_p))$ 
  // Update parameters according to gradients
   $\theta_{Enc} \leftarrow -\nabla_{\theta_{Enc}} (L_{prior} + L_{Dis^{Dis}})$ 
   $\theta_{Dec} \leftarrow -\nabla_{\theta_{Dec}} (\gamma L_{Dis^{Dis}} - L_{GAN})$ 
   $\theta_{Dis} \leftarrow -\nabla_{\theta_{Dis}} L_{GAN}$ 
until deadline

```

(a) Additional paper algorithm pseudo code

Figure 3: Additional paper algorithm



(b) Additional paper algorithm scheme

3 Data

We made experiments for MNIST and CIFAR-10 datasets.

4 Problems and solutions

Problems

- Different formulae in the text and the pseudo code in the original paper
- 1 epoch even for MNIST takes about 17 minutes
- We had no time to try supervised setting

Solutions

- Using other papers and do updates according to the our own understanding
- Several updates of encoder-decoder
- Adjusting learning rates

5 Results

5.1 Results for AAE, Alexander

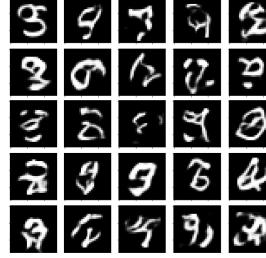


Figure 4: Samples from AAE on MNIST

5.2 Results for original paper

5.3 Alexander's results

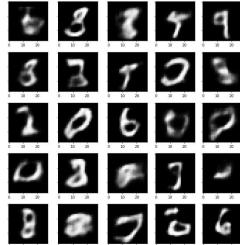
Training approach

- AAE-phase:

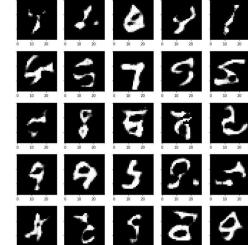
- Encoder update: $-\log(D_{code}(fake) + \varepsilon)) + L2$
- Code disc update (every 5 steps):
 - $-(\log(D_{code}(real) + \varepsilon) + \log(1 - D_{code}(fake) + \varepsilon))$
- Decoder update: $L2$

- Prior improvement phase (once in epoch):

- Decoder update: $-\log(D_{code}(sampled) + \varepsilon))$
- Code gen update: $-\log(D_{code}(sampled) + \varepsilon))$
- Image disc update: $-(\log(D_{code}(real) + \varepsilon) + \log(1 - D_{code}(fake) + \varepsilon) + \log(1 - D_{code}(sampled) + \varepsilon))$



(a) MNIST, 8 epochs



(b) MNIST, 100 epochs

Figure 5: Original paper algorithm results for different epochs

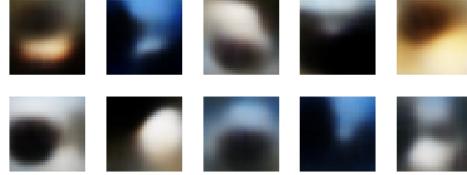


Figure 6: Results for CIFAR for original paper approach

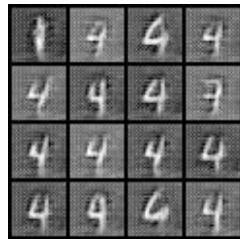


Figure 7: Samples of autoencoder work. Left image is real, right is reconstructed.

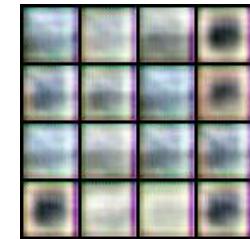
5.4 Natalia's results

Training approach can be seen in git.

The most enchanting thing here is evolution with epochs. The results for CIFAR are not good at all. Here the latent dimension is 8 the same as for MNIST.



(a) MNIST, 67 epochs



(b) CIFAR, 55 epochs

Figure 8: Original paper algorithm results for the latest epochs

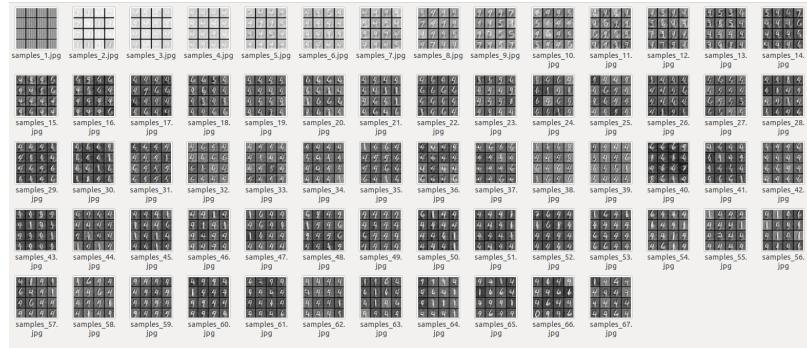


Figure 9: Results evolution for original paper. MNIST

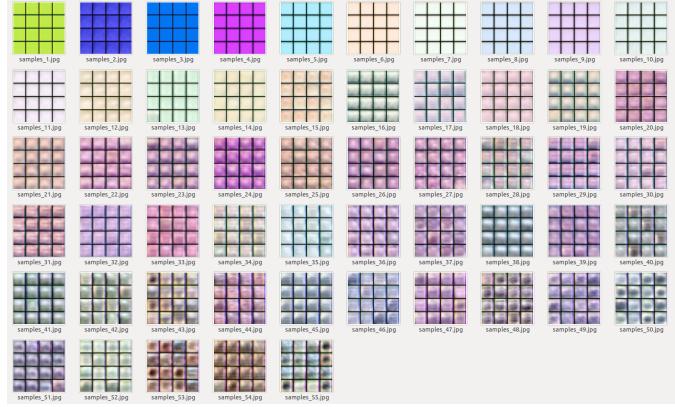


Figure 10: Results evolution for original paper. CIFAR

5.5 Results for additional algorithm

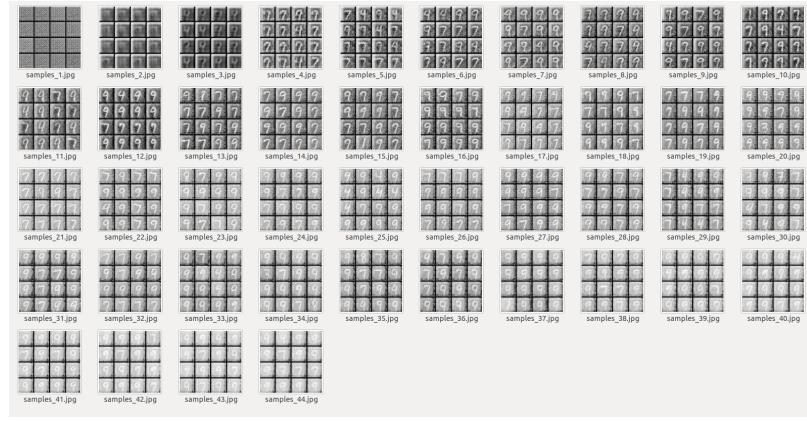


Figure 11: Results evolution for additional algorithm. MNIST

5.6 Conclusions

- The adversarial training is difficult
- There are a lot of different schemes of generative models combining GAN and VAE

5.7 Contributions

- Belozerova Polina: reading papers, overview of mathematical model of VAE and GAN for structure prediction
- Safin Alexander: reading papers, implement the classical AAE, experiments for MNIST for classical AAE, MNIST and CIFAR for original paper approach, corresponding parts in the presentation and report
- Pavlovskaya Natalia: reading papers, implement the original paper, implement the additional algorithm, experiments for MNIST for both algorithms, experiments for CIFAR-10 for original paper, corresponding parts in the presentation and report

References

- [1] Hui-Po Wang. Learning priors for adversarial autoencoders. 2018.
- [2] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, and Ian J. Goodfellow. Adversarial autoencoders. *CoRR*, abs/1511.05644, 2015.
- [3] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, and Ole Winther. Autoencoding beyond pixels using a learned similarity metric. *CoRR*, abs/1512.09300, 2015.