Deep generative modeling: Implicit models

Jakub M. Tomczak Deep Learning



#### **TYPES OF GENERATIVE MODELS**



	Training	Likelihood	Sampling	Compression
Autoregressive models (e.g., PixelCNN)	Stable	Exact	Slow	Νο
Flow-based models (e.g., RealNVP)	Stable	Exact	Fast/Slow	Νο
Implicit models (e.g., GANs)	Unstable	Νο	Fast	Νο
Prescribed models (e.g., VAEs)	Stable	Approximate	Fast	Yes



	Training	Likelihood	Sampling	Compression
Autoregressive models (e.g., PixelCNN)	Stable	Exact	Slow	Νο
Flow-based models (e.g., RealNVP)	Stable	Exact	Fast/Slow	Νο
Implicit models (e.g., GANs)	Unstable	Νο	Fast	Νο
Prescribed models (e.g., VAEs)	Stable	Approximate	Fast	Yes



#### **DENSITY NETWORKS**



1. 
$$\mathbf{z} \sim p_{\lambda}(\mathbf{z})$$
  
2.  $\mathbf{x} \sim p_{\theta}(\mathbf{x} \mid \mathbf{z})$ 

The log-likelihood function:

$$\log p(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x} \mid \mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}$$



1. 
$$\mathbf{z} \sim p_{\lambda}(\mathbf{z})$$
  
2.  $\mathbf{x} \sim p_{\theta}(\mathbf{x} \mid \mathbf{z})$ 

The log-likelihood function:

$$\log p(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x} \mid \mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}$$
$$\approx \log \frac{1}{S} \sum_{s=1}^{S} \exp \left( \log p_{\theta} \left( \mathbf{x} \mid \mathbf{z}_{s} \right) \right)$$

It could be estimated by MC samples.

7



1. 
$$\mathbf{z} \sim p_{\lambda}(\mathbf{z})$$
  
2.  $\mathbf{x} \sim p_{\theta}(\mathbf{x} \mid \mathbf{z})$   
The log-likelihood function:  
 $\log p(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x} \mid \mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}$   
 $\approx \log \frac{1}{S} \sum_{i=1}^{S} \exp \left( \log p_{\theta}(\mathbf{x} \mid \mathbf{z}_{s}) \right)$ 

It could be estimated s=1by MC samples. If we take standard Gaussian prior, <sup>8</sup> we need to model p(xlz) only!









## It must be a powerful (=flexible) transformation!



11



## It must be a powerful (=flexible) transformation! NEURAL NETWORK





# Neural network outputs parameters of a distribution, e.g., a mixture of Gaussians.



The log-likelihood function:

$$\log p(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x} \mid \mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}$$
$$\approx \log \frac{1}{S} \sum_{s=1}^{S} \exp \left( \log p_{\theta} \left( \mathbf{x} \mid \mathbf{z}_{s} \right) \right)$$



Training procedure:

- 1. Sample multiple z's from the prior (e.g., standard Gaussian).
- 2. Apply log-sum-exp-trick, and apply backpropagation.



The log-likelihood function:

$$\log p(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x} \mid \mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}$$
$$\approx \log \frac{1}{S} \sum_{s=1}^{S} \exp \left( \log p_{\theta} \left( \mathbf{x} \mid \mathbf{z}_{s} \right) \right)$$



Training procedure:

- 1. Sample multiple z's from the prior (e.g., standard Gaussian).
- 2. Apply log-sum-exp-trick, and apply backpropagation.

Drawback: It scales badly in high-dimensional cases...



#### **Advantages**

- ✓ Non-linear transformations.
- ✓ Allows to generate.

#### **Disadvantages**

- No analytical solutions.
- No exact likelihood.
- It requires a lot of samples from the prior.
- Fails in high-dim.
- It requires an explicit distribution (e.g., Gaussian).



#### **Advantages**

- ✓ Non-linear transformations.
- ✓ Allows to generate.

#### **Disadvantages**

- No analytical solutions.
- No exact likelihood.
- It requires a lot of samples from the prior.
- Fails in high-dim.
- It requires an explicit distribution
  - (e.g., Gaussian).

Can we do better?

16

#### IMPLICIT DISTRIBUTIONS



Let us look again at the Density Network model.

The idea is to inject noise to a neural network that serves as a generator:





Let us look again at the Density Network model.

The idea is to inject noise to a neural network that serves as a generator:



But now, we don't specify the distribution (e.g., MoG), but use a flexible transformation directly to return an image. This is now implicit. VU



It defines an **implicit distribution** (i.e., we do not assume any form of it), and it could be seen as Dirac's delta:

$$p(\mathbf{x} \,|\, \mathbf{z}) = \delta\left(\mathbf{x} - f(\mathbf{z})\right)$$



It defines an **implicit distribution** (i.e., we do not assume any form of it), and it could be seen as Dirac's delta:

$$p(\mathbf{x} \mid \mathbf{z}) = \delta\left(\mathbf{x} - f(\mathbf{z})\right)$$

However, now we cannot use the likelihood-based approach, because  $\ln \delta (\mathbf{x} - f(\mathbf{z}))$  is ill-defined.



It defines an **implicit distribution** (i.e., we do not assume any form of it), and it could be seen as Dirac's delta:

$$p(\mathbf{x} \mid \mathbf{z}) = \delta\left(\mathbf{x} - f(\mathbf{z})\right)$$

However, now we cannot use the likelihood-based approach, because  $\ln \delta (\mathbf{x} - f(\mathbf{z}))$  is ill-defined.

We need to use a different approach.



#### GENERATIVE ADVERSARIAL NETWORKS













An art expert





#### ... and a real artist





An art expert



Let's imagine two actors:



The fraud aims to copy the real artist and cheat the art expert.



An art expert



... and a real artist



Let's imagine two actors:



The fraud aims to copy the real artist and cheat the art expert.



An art expert

The expert assesses a painting and gives her opinion.



#### ... and a real artist



Let's imagine two actors:



The fraud aims to copy the real artist and cheat the art expert.

The fraud learns and tries to fool the expert.

A fraud



... and a real artist



An art expert

The expert assesses a painting and gives her opinion.































1. Sample **z**.

2. Generate *G*(**z**).

3. Discriminate whether given image is real or fake.







1. Sample **z**.

2. Generate *G*(**z**).

3. Discriminate whether given image is real or fake.

What about the learning objective?





 $\min_{G} \max_{D} \mathbb{E}_{\mathbf{x} \sim p_{real}}[\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))]$ 



$$\min_{G} \max_{D} \mathbb{E}_{\mathbf{x} \sim p_{real}}[\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))]$$

# It resemblances the logarithm of the Bernoulli distribution: $y \log p(y = 1) + (1 - y)\log(1 - p(y = 1))$



$$\min_{G} \max_{D} \mathbb{E}_{\mathbf{x} \sim p_{real}}[\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))]$$
  
It resemblances the logarithm of the Bernoulli distribution:  
$$y \log p(y = 1) + (1 - y) \log(1 - p(y = 1))$$



$$\begin{split} \min_{G} \max_{D} \mathbb{E}_{\mathbf{x} \sim p_{real}}[\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))] \\ \end{split}$$
It resemblances the logarithm of the Bernoulli distribution:  $y \log p(y = 1) + (1 - y) \log(1 - p(y = 1)) \end{split}$ 

Therefore, the discriminator network should end with a sigmoid function to mimic probability.



$$\min_{G} \max_{D} \mathbb{E}_{\mathbf{x} \sim p_{real}}[\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))]$$

We want to minimize wrt. generator.



We want to maximize wrt. discriminator.



- 1. Sample **z**.
- 2. Generate *G*(**z**).





- 1. Sample z.
- 2. Generate *G*(**z**).

The learning objective (adversarial loss):

 $\min_{G} \max_{D} \mathbb{E}_{\mathbf{x} \sim p_{real}}[\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))]$ 





- 1. Sample z.
- 2. Generate *G*(**z**).

The learning objective (adversarial loss):

$$\min_{G} \max_{D} \mathbb{E}_{\mathbf{x} \sim p_{real}}[\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))]$$

#### Learning:

1. Generate fake images, and minimize wrt. G.

2. Take real & fake images, and maximize wrt. D.



import torch.nn as nn

```
class GAN(nn.Module):
    def __init__(self, D, M):
        super(GAN, self).__init__()
        self.D = D
        self.M = M
```

```
self.gen1 = nn.Linear(in_features= self.M, out_features=300)
self.gen2 = nn.Linear(in_features=300, out_features= self.D)
```

```
self.dis1 = nn.Linear(in_features= self.D, out_features=300)
self.dis2 = nn.Linear(in_features=300, out_features=1)
```



```
def generate(self, N):
    z = torch.randn(size=(N, self.D))
    x_gen = self.gen1(z)
    x_gen = nn.functional.relu(x_gen)
    x_gen = self.gen2(x_gen)
    return x_gen
```

```
def discriminate(self, x):
    y = self.dis1(x)
    y = nn.functional.relu(y)
    y = self.dis2(y)
    y = torch.sigmoid(y)
    return y
```



```
def gen loss(self, d gen):
   return torch.log(1. - d gen)
def dis loss(self, d real, d gen):
   # We maximize wrt. the discriminator, but optimizers minimize!
   # We need to include the negative sign!
   return -(torch.log(d real) + torch.log(1. - d gen))
def forward(self, x real):
   x gen = self.generate(N=x real.shape[0])
   d real = self.discriminate(x real)
   d gen = self.discriminate(x gen)
```

```
return d_real, d_gen
```



```
def gen loss(self, d gen):
   return torch.log(1. - d gem)
def dis loss(self, d real, d gen):
   # We maximize wrt. the discriminator, but optimizers minimize!
   # We need to include the negative sign!
   return -(torch.log(d real) + torch.log(1. - d gen))
def forward(self, x real):
   x gen = self.generate(N=x real.shape[0])
   d real = self.discriminate(x real)
   d gen = self.discriminate(x gen)
   return d real, d gen
```

We can use two optimizers, one for d\_real & d\_gen, and one for d\_gen.

#### **GENERATIONS**



Training Data

Samples



Salimans, T., et al. (2016). Improved techniques for training GANs. NeurIPS

53

#### **Advantages**

- ✓ Non-linear transformations.
- ✓ Allows to generate.
- ✓ Learnable loss.
- ✓ Allows implicit models.
- ✓ Works in high-dim.

#### **Disadvantages**

- No exact likelihood.
- Unstable training.
- Missing mode problem (i.e., it doesn't cover the whole space).
- No clear way for quantitative assessment.





For instance, we can use the **earth-mover distance**:

$$\min_{G} \max_{D \in \mathcal{W}} \mathbb{E}_{\mathbf{x} \sim p_{real}}[D(\mathbf{x})] - \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[D(G(\mathbf{z}))]$$

where the discriminator is a 1-Lipschitz function.



Arjovsky, M., Chintala, S., & Bottou, L. (2017). Wasserstein GAN. arXiv preprint arXiv:1701.07875.

For instance, we can use the **earth-mover distance**:

$$\min_{G} \max_{D \in \mathcal{W}} \mathbb{E}_{\mathbf{x} \sim p_{real}}[D(\mathbf{x})] - \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[D(G(\mathbf{z}))]$$

where the discriminator is a 1-Lipschitz function.

We need to clip weights of the discriminator: clip(weights, -c, c)



Arjovsky, M., Chintala, S., & Bottou, L. (2017). Wasserstein GAN. arXiv preprint arXiv:1701.07875.

For instance, we can use the **earth-mover distance**:

$$\min_{G} \max_{D \in \mathcal{W}} \mathbb{E}_{\mathbf{x} \sim p_{real}}[D(\mathbf{x})] - \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[D(G(\mathbf{z}))]$$

where the discriminator is a 1-Lipschitz function.

We need to clip weights of the discriminator: clip(weights, -c, c)

## It stabilizes training, but other problems remain.



Arjovsky, M., Chintala, S., & Bottou, L. (2017). Wasserstein GAN. arXiv preprint arXiv:1701.07875.



