



## Lecture Notes - GANs

Muskan - 190520

Videeta Sharma - 190959

### Neural Network

First, we were introduced to the structure of a simple neural network. The building block of Neural Networks are neurons. A neuron is an entity that takes inputs, performs mathematical operations on them and returns an output. A simple neural network for performing regression with one hidden layer can be depicted as:

$$h_{i_1} = \sum_{i_0} (w_{i_1 i_0} x_{i_0}) + b_{i_1} \quad (1)$$

$$v_{i_1} = \sigma(h_{i_1}) \quad (2)$$

$$y_{i_3} = \sum_{i_2} (w_{i_3 i_2} v_{i_2}) + b_{i_3} \quad (3)$$

where  $\sigma$  is the sigmoid function, given as:

$$\sigma(h) = \frac{1}{1 + \exp(-h)} \quad (4)$$

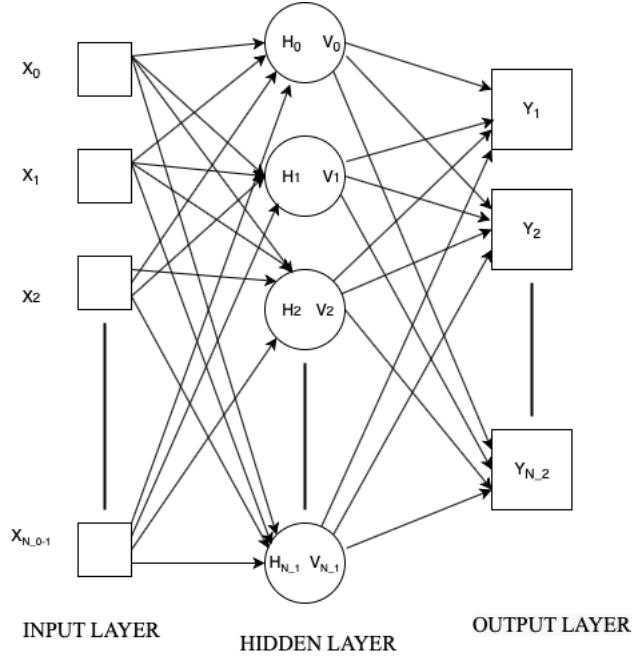
$x_{i_0}$  are the inputs to the neural network, where  $0 \leq i_0 < N_0$ . These inputs are fed to the hidden layer where the neurons give the output  $h_{i_1}$ ,  $1 \leq i_1 \leq N_1$ .

$b_{i_1}$  is the bias term associated with the hidden layer and  $w_{i_1 i_0}$  are the weights associated with neurons in the hidden layer. To introduce non-linearity, sigmoid function is used. The final output is calculated by using  $h_{i_1}$  as inputs,  $w_{i_3 i_2}$  as weights and  $b_{i_3}$  as the bias terms associated with the output layer.

The trainable parameters are represented by  $\theta = [w, b]$ . Taking Mean Square error as the Loss function for this model,

$$L(\theta) = \sum_{data} \sum_{i_3} (y_{i_3}^d - y_{i_3})^2 \quad (5)$$

$$L(\theta) = \mathbb{E}_{(x, y_d) \sim data} \left[ \sum_{i_3} (y_{i_3}^d - y_{i_3})^2 \right] \quad (6)$$



### **STRUCTURE OF A SIMPLE NEURAL NETWORK WITH ONE HIDDEN LAYER**

Figure 1: Depiction of Neural Network with one hidden layer.

If we want to use neural networks for Classification, the sigmoid function can be used in the last layer, instead of (3).

$$y_{i3} = \sigma(v_{i3}) \quad (7)$$

For performing multilabel classification, Binary Cross Entropy Loss function is used.

$$L(\theta) = - \sum_{data} \sum_{i3} (y_{i3}^d \log(y_{i3}) + (1 - y_{i3}^d) \log(1 - y_{i3})) \quad (8)$$

$$L(\theta) = - \mathbb{E}_{(x,y_{i3}) \sim data} \left[ \sum_{i3} (y_{i3}^d \log(y_{i3}) + (1 - y_{i3}^d) \log(1 - y_{i3})) \right] \quad (9)$$

$$L(\theta) = - \sum_{i3} \left( \mathbb{E}_{(x,y) \sim data, y_{i3}^d = 1} \log(y_{i3}) + \mathbb{E}_{(x,y) \sim data, y_{i3}^d = 0} (1 - \log(y_{i3})) \right) \quad (10)$$

### **Generative Adversarial Networks**

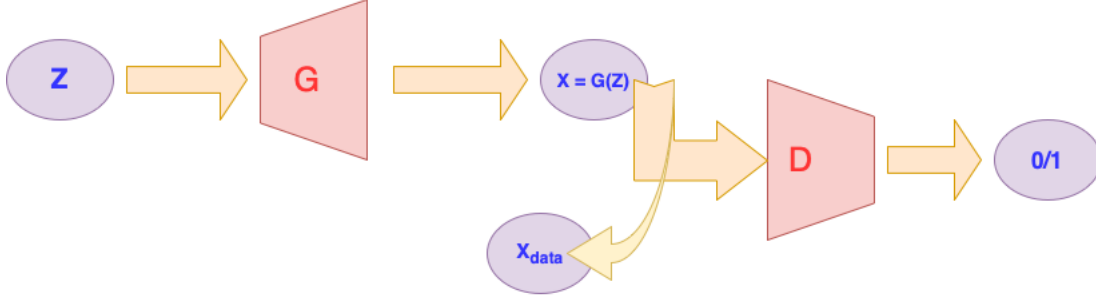
If we are given samples from a Complicated Probability Distribution  $p(x)$ , then to sample more values from unknown  $p(x)$ , we can use Generative Adversarial Networks (GANs).

It consists of a Generator and Discriminator. A generator takes samples  $z$  from a simple, known probability distribution  $p(z)$  and generates  $x$ .

$$\vec{x} = G(\vec{z}) \quad (11)$$

Since  $G$  is not a bijective function, we cannot calculate the probability distribution of  $\vec{x}$  by calculating Jacobian of the inverse transform. Consider  $\theta_g$  as the trainable parameters of the Generator.

Now, we connect a Discriminator which takes  $\vec{x}$  and  $\vec{x}_{data}$  as the input and gives output 0 and 1 for them respectively. The discriminator works as the classifier that differentiates between real and false data.



## STRUCTURE OF A GAN

Consider  $\theta_d$  as the trainable parameters of the Discriminator represented by  $D$ . The loss function is given by:

$$L(\theta_g, \theta_d) = -\left[ \mathbb{E}_{(x, y^d) \sim \text{data}, y^d=1, \vec{x}=\vec{x}_{data}} [\log(D(\vec{x}))] + \mathbb{E}_{\vec{x}=G(\vec{z}), \vec{z} \sim p(\vec{z})} [\log(1 - D(\vec{x}))] \right] \quad (12)$$

On simplifying this, we can write

$$L(\theta_g, \theta_d) = -\left[ \mathbb{E}_{\vec{x} \sim p_{data}(\vec{x})} [\log(D(\vec{x}))] + \mathbb{E}_{\vec{z} \sim p_z(\vec{z})} [\log(1 - D(G(\vec{z})))] \right] \quad (13)$$

For training, the Loss Function has to be minimised with respect to  $\theta_d$ , which enables the Discriminator to classify accurately, and maximised with respect to  $\theta_g$  which enables the Generator to give more accurate, real-like samples. Between these 2 pursuits, maximising  $\theta_g$  is treated as the priority.

## Training Algorithm for GANs

- Run a loop for  $N$  iterations and perform the following steps:
  - Perform  $K$  iterations of the following:
    - Sample a minibatch of  $m$  samples,  $\vec{z} \sim p_z(\vec{z})$ ,  $\vec{x} = G(\vec{z})$ .
    - Sample another minibatch of  $m$  samples of  $\vec{x}_{data} \sim p_{data}(\vec{x})$
    - Update parameters of Discriminator  $\theta_d = \theta_d - \eta \nabla_{\theta_d} L(\theta_g, \theta_d)$ .
  - Sample a minibatch of  $m$  samples of  $\vec{z} \sim p_z(\vec{z})$ .
  - Update parameters of Generator  $\theta_g = \theta_g + \eta \nabla_{\theta_g} L(\theta_g, \theta_d)$ . For the generator's parameters, the second term of loss function will be considered i.e  $\nabla_{\theta_g} \sum_{s=1}^m \log(1 - D(G(\vec{z}^s)))$ . The paper suggests that instead of this, we can use  $-\nabla_{\theta_g} \sum_{s=1}^m \log(D(G(\vec{z}^s)))$  for simplifying the calculation. Therefore, we can perform the following update operation,  $\theta_g = \theta_g - \eta \nabla_{\theta_g} \sum_{s=1}^m \log(D(G(\vec{z}^s)))$ .

Here  $\eta$  is the learning rate, which is a hyperparameter and retains a constant value during the training of the model. It is optimised by testing for different values on the validation data set.

The update equations for trainable parameters are used to steer the parameters towards their optimised value by using the partial derivatives.

## Theoretical Analysis

- The minimax game (13) has a global optimum for  $p_g = p_{data}$ .
- The above algorithm obtains the desired result.

1. For G fixed, the optimal discriminator D is

$$D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)} \quad (14)$$

**Proof:**

$$L(\theta_g, \theta_d) = -[\mathbb{E}_{\vec{x} \sim p_{data}(\vec{x})} [\log(D(\vec{x}))] + \mathbb{E}_{\vec{x} \sim p_g(\vec{x})} [\log(1 - D(\vec{x}))]] \quad (15)$$

$$L(\theta_g, \theta_d) = - \int (p_{data}(x) \log(D(x)) + p_g(x) \log(1 - D(x))) dx \quad (16)$$

We want to minimize the function with respect to D. The function is of the form  $f(y) = -(a \cdot \log(y) + b \cdot \log(1 - y))$ . To find minimum value, take the derivative. The function will be minimum at

$$y = \frac{a}{a + b} \quad (17)$$

Using (16) and (17),

$$D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)} \quad (18)$$

2. For fixed  $D = D_G^*(x)$ , optimal G is

$$\min_{\theta_g} L(\theta_g, \theta_d) = \mathbb{E}_{\vec{x} \sim p_{data}(\vec{x})} [\log(D_G^*(x))] + \mathbb{E}_{\vec{z} \sim p_g(\vec{z})} [\log(1 - D_G^*(x))] \quad (19)$$

$$L(\theta_g, \theta_d) = \mathbb{E}_{\vec{x} \sim p_{data}(\vec{x})} [\log(\frac{p_{data}(x)}{p_{data}(x) + p_g(x)})] + \mathbb{E}_{\vec{z} \sim p_g(\vec{z})} [\log(\frac{p_g(x)}{p_{data}(x) + p_g(x)})] \quad (20)$$

As  $p_g$  is the function of generator, we want to optimize the above equation with respect to it. To optimize, a trick is used. The first term looks like KL divergence. The only problem is that the denominator is not a pdf. To convert it into a pdf, we can take its average. So, inside the log term, divide both numerator and denominator by 2.

$$L(\theta_g, \theta_d) = \mathbb{E}_{\vec{x} \sim p_{data}(\vec{x})} [\log(\frac{\frac{p_{data}(x)}{2}}{\frac{p_{data}(x) + p_g(x)}{2}})] + \mathbb{E}_{\vec{z} \sim p_g(\vec{z})} [\log(\frac{\frac{p_g(x)}{2}}{\frac{p_{data}(x) + p_g(x)}{2}})] \quad (21)$$

$$L(\theta_g, \theta_d) = -\log 2 + KL[p_{data}(x) \parallel \frac{p_{data}(x) + p_g(x)}{2}] - \log 2 + KL[p_g(x) \parallel \frac{p_{data}(x) + p_g(x)}{2}] \quad (22)$$

Minimum value of KL divergence is zero when both the pdfs are equal. The sum of these two KL divergence's is called **Jensen-Shannon divergence**. It can be said that JSD is a symmetric variant of KL Divergence. So now (22) is transformed to

$$L(\theta_g, \theta_d) = -2\log 2 + 2JS D(p_{data} \parallel p_g) \quad (23)$$

The value of JSD is also zero when both the pdfs are equal. So value of JSD is zero when  $p_{data} = p_g$ . This is the condition for **global optimum**.

## Limitations of GANs

Consider the case when  $G(z)$  is mapped to a valid  $x \in \text{data}$  for all  $z$ . This is called **Mode Collapse**.

$$L(\theta_g, \theta_d) = -\left[ \mathbb{E}_{\vec{x} \sim p_{data}(\vec{x})} [\log(D_G(x))] + \mathbb{E}_{\vec{z} \sim p_z(\vec{z})} [\log(1 - D_G(x))] \right] \quad (24)$$

As  $G(z) = x_1$  value of  $D(x)$  will be 1. So first term in the loss function will be zero and second term is  $-\infty$ . So because of negative sign, loss becomes very large.

## Reference

- Generative Adversarial Nets: Ian J. Goodfellow