

Using Deep Learning To Implement Domain Generation Algorithm

Student: Chia-Ruei Lee

Advisor: Chia-Mu Yu



CONTENTS

- A. Introduction
 - 1. Botnet
 - 2. Domain Generation Algorithm
- B. Deep Learning
 - 1. Neural Network
 - 2. Auto-Encoder
 - 3. Generative Adversarial Networks (GAN)
 - 4. Wasserstein GAN (WGAN)
 - 5. Improved Wasserstein GAN
- C. Using GAN To Implement DGA
 - 1. Requirement



Botnet

Introduction



Domain Generation Algorithm

Introduction

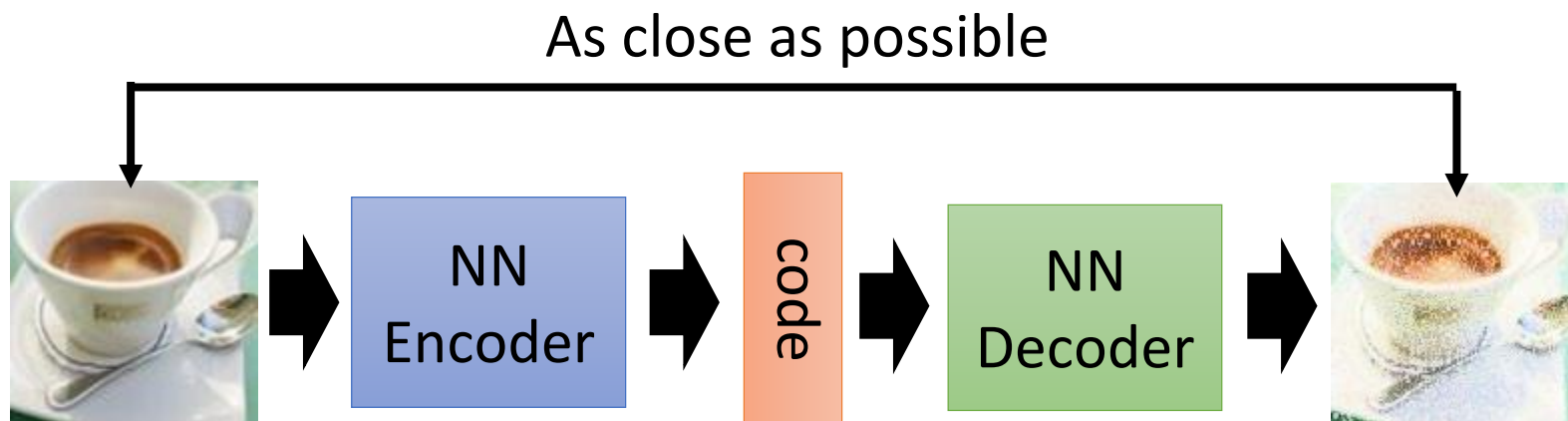
Neural Network

Deep Learning

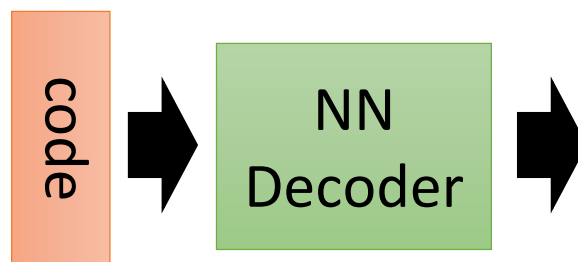
6 Auto-Encoder

Deep Learning

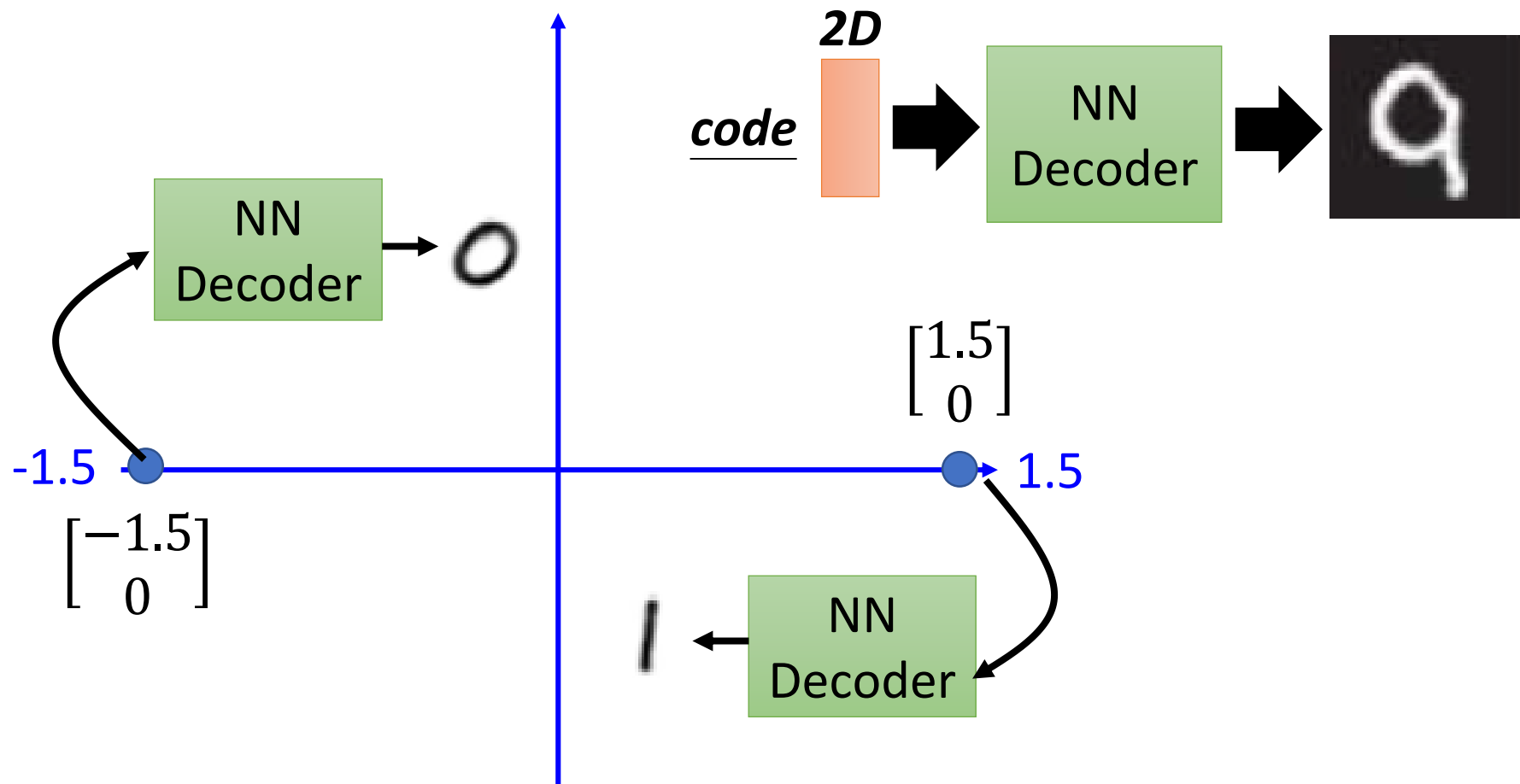
AUTO-ENCODER



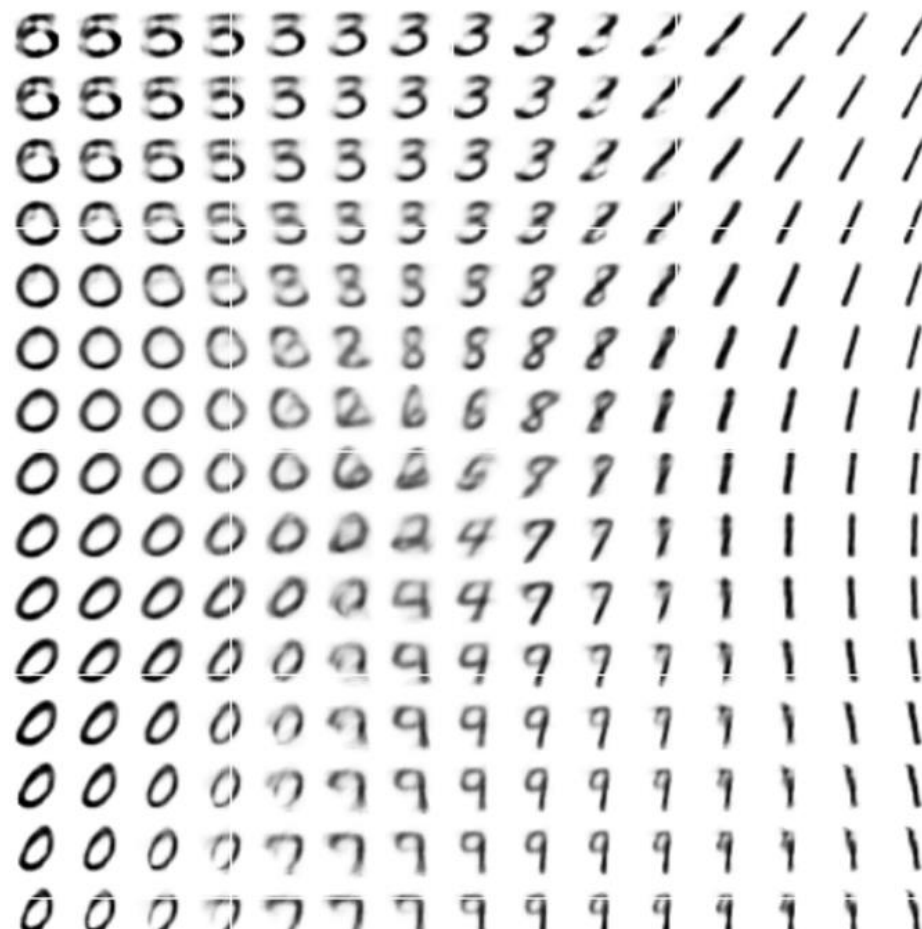
Randomly generate
a vector as code



AUTO-ENCODER



AUTO-ENCODER



From Hung-yi Lee, “Generative Adversarial Network/Basic Idea/Machine Learning and having it deep and structured (2017, Spring)”
http://speech.ee.ntu.edu.tw/~tlkagk/courses_MLDS17.html



Generative Adversarial Networks (GAN)

10

Deep Learning

ALGORITHM

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of m examples $\{x^{(1)}, \dots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(x^{(i)}) + \log (1 - D(G(z^{(i)}))) \right].$$

end for

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^{(i)}))).$$

end for

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

12

Wasserstein GAN (WGAN)

Deep Learning

ALGORITHM

Algorithm 1 WGAN, our proposed algorithm. All experiments in the paper used the default values $\alpha = 0.00005$, $c = 0.01$, $m = 64$, $n_{\text{critic}} = 5$.

Require: : α , the learning rate. c , the clipping parameter. m , the batch size. n_{critic} , the number of iterations of the critic per generator iteration.

Require: : w_0 , initial critic parameters. θ_0 , initial generator's parameters.

```
1: while  $\theta$  has not converged do
2:   for  $t = 0, \dots, n_{\text{critic}}$  do
3:     Sample  $\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$  a batch from the real data.
4:     Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
5:      $g_w \leftarrow \nabla_w [\frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))]$ 
6:      $w \leftarrow w + \alpha \cdot \text{RMSProp}(w, g_w)$ 
7:      $w \leftarrow \text{clip}(w, -c, c)$ 
8:   end for
9:   Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
10:   $g_\theta \leftarrow -\nabla_\theta \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))$ 
11:   $\theta \leftarrow \theta - \alpha \cdot \text{RMSProp}(\theta, g_\theta)$ 
12: end while
```

From “Wasserstein GAN”

<https://arxiv.org/abs/1701.07875>

14

Improved Wasserstein GAN (WGAN-GP)

Deep Learning

ALGORITHM

Algorithm 1 WGAN with gradient penalty. We use default values of $\lambda = 10$, $n_{\text{critic}} = 5$, $\alpha = 0.0001$, $\beta_1 = 0$, $\beta_2 = 0.9$.

Require: The gradient penalty coefficient λ , the number of critic iterations per generator iteration n_{critic} , the batch size m , Adam hyperparameters α, β_1, β_2 .

Require: initial critic parameters w_0 , initial generator parameters θ_0 .

```
1: while  $\theta$  has not converged do
2:   for  $t = 1, \dots, n_{\text{critic}}$  do
3:     for  $i = 1, \dots, m$  do
4:       Sample real data  $\mathbf{x} \sim \mathbb{P}_r$ , latent variable  $\mathbf{z} \sim p(\mathbf{z})$ , a random number  $\epsilon \sim U[0, 1]$ .
5:        $\tilde{\mathbf{x}} \leftarrow G_{\theta}(\mathbf{z})$ 
6:        $\hat{\mathbf{x}} \leftarrow \epsilon \mathbf{x} + (1 - \epsilon) \tilde{\mathbf{x}}$ 
7:        $L^{(i)} \leftarrow D_w(\tilde{\mathbf{x}}) - D_w(\mathbf{x}) + \lambda(\|\nabla_{\hat{\mathbf{x}}} D_w(\hat{\mathbf{x}})\|_2 - 1)^2$ 
8:     end for
9:      $w \leftarrow \text{Adam}(\nabla_w \frac{1}{m} \sum_{i=1}^m L^{(i)}, w, \alpha, \beta_1, \beta_2)$ 
10:   end for
11:   Sample a batch of latent variables  $\{\mathbf{z}^{(i)}\}_{i=1}^m \sim p(\mathbf{z})$ .
12:    $\theta \leftarrow \text{Adam}(\nabla_{\theta} \frac{1}{m} \sum_{i=1}^m -D_w(G_{\theta}(\mathbf{z}^{(i)})), \theta, \alpha, \beta_1, \beta_2)$ 
13: end while
```

From “Improved Training of Wasserstein GANs”
<https://arxiv.org/abs/1704.00028>



Requirement

Using GAN To Implement DGA

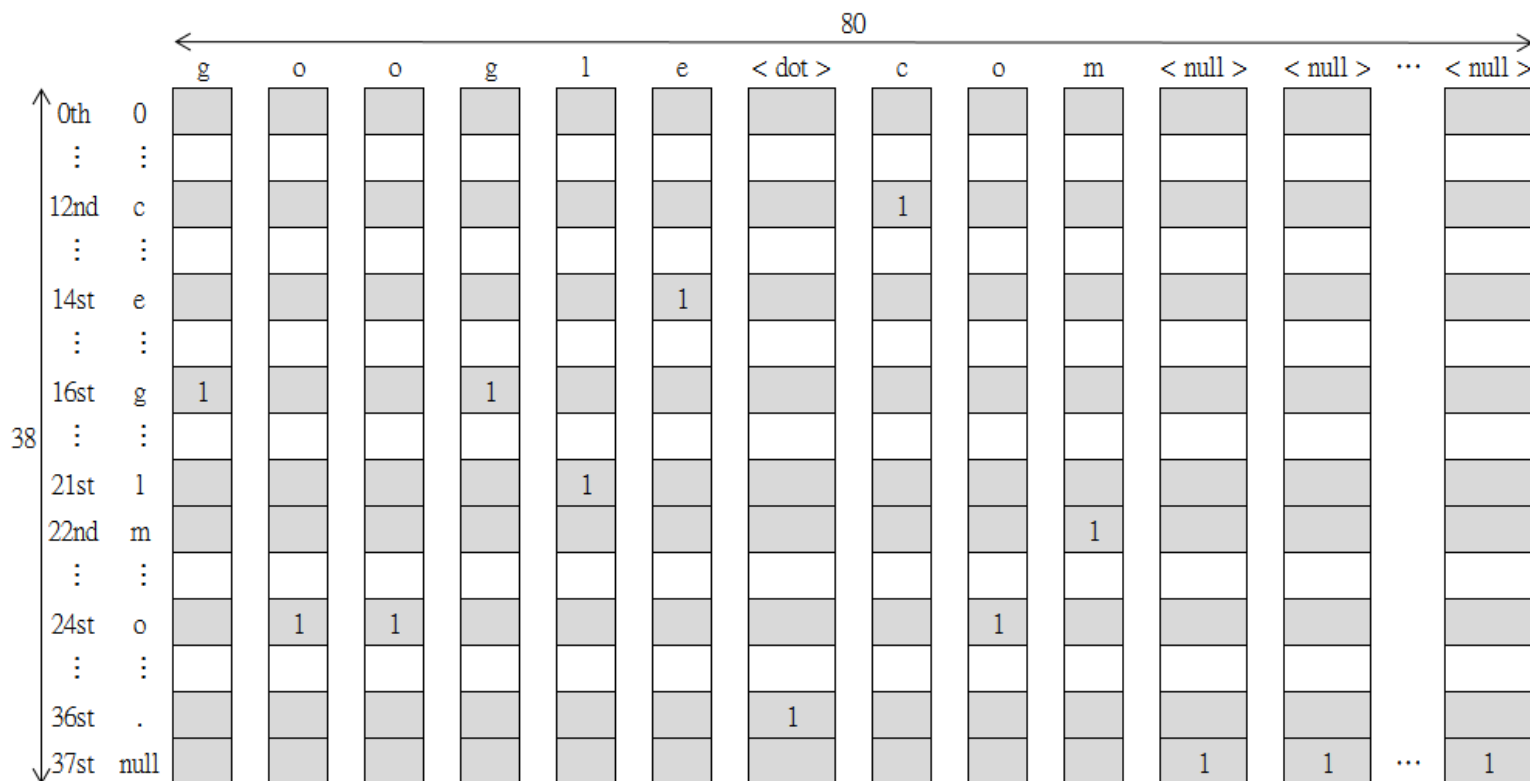
REQUIREMENT

- Q: Where is my training data?
- A: From Alexa Top 1-million sites.

	A	B
1	1	google.com
2	2	youtube.com
3	3	facebook.com
4	4	baidu.com
5	5	yahoo.com
6	6	wikipedia.org
7	7	amazon.com
8	8	twitter.com
9	9	qq.com
10	10	google.co.in
11	11	live.com
12	12	taobao.com
13	13	bing.com
14	14	google.co.jp
15	15	msn.com
16	16	yahoo.co.jp
17	17	linkedin.com
18	18	sina.com.cn
19	19	weibo.com
20	20	vk.com
21	21	instagram.com
22	22	google.ru
23	23	yandex.ru
24	24	google.de
25	25	hao123.com
26	26	ebay.com
27	27	reddit.com
28	28	google.co.uk
29	29	amazon.co.jp
30	30	t.co
31	31	google.com.br
32	32	mail.ru
33	33	google.fr
34	34	pinterest.com
35	35	netflix.com

REQUIREMENT

- Q: How to use this data?
- A: I deal these domain names as 38x80 pictures.



REQUIREMENT

- Because we don't have the "code". We have to train an Auto-encoder model.

