## Applied Deep Learning



## Guest Lecture by Hung-yi Lee Generative Adversarial Network



May 19th, 2020 http://adl.miulab.tw



#### **Three Categories of GAN**



## Generative Adversarial Network (GAN)

• Anime face generation as example



## Algorithm

• Initialize generator and discriminator

 $\square$ 

- In each training iteration:
- **Step 1**: Fix generator G, and update discriminator D



Discriminator learns to assign high scores to real objects and low scores to generated objects.

## Algorithm

Initialize generator and discriminator

G

• In each training iteration:

**Step 2**: Fix discriminator D, and update generator G

Generator learns to "fool" the discriminator



## Algorithm

• Initialize generator and discriminator

G

D

• In each training iteration:





https://crypko.ai/#/

## GAN is hard to train .....



(I found this joke from 陳柏文's facebook.)

## Three Categories of GAN



unpaired data

## Text-to-Image



a bird is flying



Traditional supervised approach



[Scott Reed, et al, ICML, 2016]

## Conditional GAN





Generator will learn to generate realistic images ....

But completely ignore the input conditions.



[Scott Reed, et al, ICML, 2016]









## Conditional GAN - Sound-to-image

The images are generated by Chia-Hung Wan and Shun-Po Chuang. https://wjohn1483.github.io/ audio\_to\_scene/index.html

# Louder Audio-to-image



## Conditional GAN - Image-to-label

#### Multi-label Image Classifier





## Conditional GAN - Image-to-label

The classifiers can have different architectures.

The classifiers are trained as conditional GAN.

[Tsai, et al., submitted to ICASSP 2019]

F1	MS-COCO	NUS-WIDE	
VGG-16	56.0	33.9	
+ GAN	60.4	41.2	
Inception	62.4	53.5	
+GAN	63.8	55.8	
Resnet-101	62.8	53.1	
+GAN	64.0	55.4	
Resnet-152	63.3	52.1	
+GAN	63.9	54.1	
Att-RNN	62.1	54.7	
RLSD	62.0	46.9	

## Conditional GAN - Image-to-label

The classifiers can have different architectures.

The classifiers are trained as conditional GAN.

Conditional GAN outperforms other models designed for multi-label.

F1	MS-COCO	NUS-WIDE	
VGG-16	56.0	33.9	
+ GAN	60.4	41.2	
Inception	62.4	53.5	
+GAN	63.8	55.8	
Resnet-101	62.8	53.1	
+GAN	64.0	55.4	
Resnet-152	63.3	52.1	
+GAN	63.9	54.1	
Att-RNN	62.1	54.7	
RLSD	62.0	46.9	

## Talking Head



https://arxiv.org/abs/1905.08233

## **Three Categories of GAN**



## Cycle GAN



#### Domain Y



#### Domain X







Input image belongs to domain Y or not

Domain Y



Domain Y

[Jun-Yan Zhu, et al., ICCV, 2017]

## Cycle GAN

#### as close as possible



Domain Y

## Cycle GAN

#### as close as possible







## Three Categories of Solutions

#### Gumbel-softmax

• [Matt J. Kusner, et al, arXiv, 2016]

#### **Continuous Input for Discriminator**

[Sai Rajeswar, et al., arXiv, 2017][Ofir Press, et al., ICML workshop, 2017][Zhen Xu, et al., EMNLP, 2017][Alex Lamb, et al., NIPS, 2016][Yizhe Zhang, et al., ICML, 2017]

#### "Reinforcement Learning"

[Yu, et al., AAAI, 2017][Li, et al., EMNLP, 2017][Tong Che, et al, arXiv, 2017][Jiaxian Guo, et al., AAAI, 2018][Kevin Lin, et al, NIPS, 2017][William Fedus, et al., ICLR, 2018]



Negative sentence to positive sentence: it's a crappy day -> it's a great day i wish you could be here -> you could be here it's not a good idea -> it's good idea i miss you -> i love you i don't love you -> i love you i can't do that -> i can do that ifeel so sad -> i happy it's a bad day -> it's a good day it's a dummy day -> it's a great day sorry for doing such a horrible thing -> thanks for doing a great thing my doggy is sick -> my doggy is my doggy my little doggy is sick -> my little doggy is my little doggy



#### Negative sentence to positive sentence:

胃疼,沒睡醒,各種不舒服-> 生日快樂,睡醒,超級舒服 我都想去上班了,真夠賤的!-> 我都想去睡了,真帥的! 暈死了,吃燒烤、竟然遇到個變態狂-> 哈哈好~,吃燒烤~竟 然遇到帥狂

我肚子痛的厲害 -> 我生日快樂厲害

感冒了,難受的說不出話來了!-> 感冒了,開心的說不出話來!





- Supervised learning needs lots of annotated speech.
- However, most of the languages are low resourced.



Listening to human talking

Reading text on the Internet

## Acoustic Token Discovery



Acoustic tokens can be discovered from audio collection without text annotation.

Acoustic tokens: chunks of acoustically similar audio segments with token IDs [Zhang & Glass, ASRU 09]

[Huijbregts, ICASSP 11] [Chan & Lee, Interspeech 11]

## Acoustic Token Discovery



Acoustic tokens can be discovered from audio collection without text annotation.

Acoustic tokens: chunks of acoustically similar audio segments with token IDs [Zhang & Glass, ASRU 09]

[Huijbregts, ICASSP 11] [Chan & Lee, Interspeech 11]

## Acoustic Token Discovery



**Phonetic-level acoustic tokens** are obtained by segmental sequence-to-sequence autoencoder.

## **Unsupervised Speech Recognition**



## Model



#### **Experimental Results**

Approaches		Matched		Nonmatched				
		(all 4000)		(3000/1000)				
		FER	PER	FER	PER			
(I) Supervised (labeled)								
(a) RNN Transducer [23]		-	17.7	-	-			
(b) standard HMMs		-	21.5	-	-			
(c) Phoneme classifier		27.0	28.9	-	-			
(II) Unsupervised (with oracle boundaries)								
(d) Relationship mapping GAN [22]		40.5	40.2	43.6	43.4			
(e) Segmental Emperical-ODM [23]		33.3	32.5	40.0	40.1			
(f) Proposed: GAN		27.6	28.5	32.7	34.3			
(III) Completely unsupervised (no label at all)								
(g) Segmental Emperical-ODM [23]		-	36.5	-	41.6			
Just constraints in the second	(h) GAN	48.3	48.6	50.3	50.0			
		(i) GAN/HMM	-	30.7	-	39.5		
	iteration 2	(j) GAN	41.0	41.0	44.3	44.3		
		(k) GAN/HMM	-	27.0	-	35.5		
	iteration 3	(I) GAN	39.7	38.4	45.0	44.2		
	iteration 3	(m) GAN/HMM	-	26.1	-	33.1		


The image is modified from: Phone recognition on the TIMIT database Lopes, C. and Perdigão, F., 2011. Speech Technologies, Vol 1, pp. 285--302.

### **Three Categories of GAN**

#### 1. Typical GAN -0.3 0.1 Generator 0.9 random vector image 2. Conditional GAN blue eyes, "Girl with red hair, Generator red hair" short hair text paired data image 3. Unsupervised Conditional GAN domain x domain y Х V Generator

Photo

unpaired data

Vincent van

Gogh's style

# Theory behind GAN

### Generation

### Using Generative Adversarial Network (GAN)



### Generation

x: an image (a highdimensional vector)

• We want to find data distribution  $P_{data}(x)$ 



### Maximum Likelihood Estimation

- Given a data distribution  $P_{data}(x)$  (We can sample from it.)
- We have a distribution  $P_G(x; \theta)$  parameterized by  $\theta$ 
  - We want to find  $\theta$  such that  $P_G(x; \theta)$  close to  $P_{data}(x)$
  - E.g.  $P_G(x; \theta)$  is a Gaussian Mixture Model,  $\theta$  are means and variances of the Gaussians

Sample  $\{x^1, x^2, ..., x^m\}$  from  $P_{data}(x)$ We can compute  $P_G(x^i; \theta)$ 

Likelihood of generating the samples

$$L = \prod_{i=1}^{m} P_G(x^i; \theta)$$

Find  $\theta^*$  maximizing the likelihood



## Maximum Likelihood Estimation = Minimize KL Divergence

$$\theta^* = \arg \max_{\theta} \prod_{i=1}^m P_G(x^i; \theta) = \arg \max_{\theta} \log \prod_{i=1}^m P_G(x^i; \theta)$$
$$= \arg \max_{\theta} \sum_{i=1}^m \log P_G(x^i; \theta) \quad \{x^1, x^2, \dots, x^m\} \text{ from } P_{data}(x)$$
$$\approx \arg \max_{\theta} E_{x \sim P_{data}}[\log P_G(x; \theta)]$$
$$= \arg \max_{\theta} \int_x P_{data}(x) \log P_G(x; \theta) dx - \int_x P_{data}(x) \log P_{data}(x) dx$$
$$= \arg \min_{\theta} KL(P_{data}||P_G) \quad \text{How to define a general } P_G?$$



x: an image (a highdimensional vector)

• A generator G is a network. The network defines a probability distribution  $P_G$ 



$$G^* = arg \min_{G} \underline{Div(P_G, P_{data})}$$
  
Divergence between distributions  $P_G$  and  $P_{data}$   
How to compute the divergence?

### Discriminator

$$G^* = \arg\min_{G} Div(P_G, P_{data})$$

Although we do not know the distributions of  $P_G$  and  $P_{data}$ , we can sample from them.



### Discriminator $G^* = \arg\min_{G} Div(P_G, P_{data})$



### Discriminator $G^* = \arg \min_{G} Div(P_G, P_{data})$



$$\max_D V(G,D)$$

$$V = E_{x \sim P_{data}}[log D(x)] + E_{x \sim P_{G}}[log(1 - D(x))]$$

• Given G, what is the optimal D\* maximizing

$$V = E_{x \sim P_{data}}[logD(x)] + E_{x \sim P_{G}}[log(1 - D(x))]$$
  
=  $\int_{x} P_{data}(x)logD(x) dx + \int_{x} P_{G}(x)log(1 - D(x)) dx$   
=  $\int_{x} [P_{data}(x)logD(x) + P_{G}(x)log(1 - D(x))] dx$   
Assume that D(x) can be any function

• Given x, the optimal D\* maximizing  $P_{data}(x)logD(x) + P_G(x)log(1 - D(x))$ 

 $\max V(G, D)$ 

$$V = E_{x \sim P_{data}}[log D(x)] + E_{x \sim P_{G}}[log(1 - D(x))]$$

- Given x, the optimal D\* maximizing  $P_{data}(x)logD(x) + P_G(x)log(1 - D(x))$ a D b D
- Find D\* maximizing: f(D) = alog(D) + blog(1 D)

$$\frac{df(D)}{dD} = a \times \frac{1}{D} + b \times \frac{1}{1 - D} \times (-1) = 0$$

$$a \times \frac{1}{D^*} = b \times \frac{1}{1 - D^*} \quad a \times (1 - D^*) = b \times D^*$$

$$a - aD^* = bD^* \quad a = (a + b)D^*$$

$$D^* = \frac{a}{a + b} \qquad \longrightarrow \qquad D^*(x) = \frac{P_{data}(x)}{P_{data}(x) + P_G(x)} < 1$$

$$\begin{split} \max_{D} V(G, D) & V(G, D) \\ \max_{D} V(G, D) &= V(G, D^{*}) \\ & = E_{x \sim P_{data}}[logD(x)] \\ & + E_{x \sim P_{G}}[log(1 - D(x))] \\ \end{split} \\ m_{D} V(G, D) &= V(G, D^{*}) \\ & D^{*}(x) = \frac{P_{data}(x)}{P_{data}(x) + P_{G}(x)} \\ & = E_{x \sim P_{data}} \left[ log \frac{P_{data}(x)}{P_{data}(x) + P_{G}(x)} \right] \\ & + E_{x \sim P_{G}} \left[ log \frac{P_{G}(x)}{P_{data}(x) + P_{G}(x)} \right] \\ & + E_{x \sim P_{G}} \left[ log \frac{P_{G}(x)}{P_{data}(x) + P_{G}(x)} \right] \\ & = \int_{x} P_{data}(x) log \frac{\frac{1}{2} P_{data}(x)}{\frac{P_{data}(x) + P_{G}(x)}{2}} dx \\ & + 2log \frac{1}{2} - 2log2 \\ & + \int_{x} P_{G}(x) log \frac{\frac{1}{2} P_{G}(x)}{\frac{P_{data}(x) + P_{G}(x)}{2}} dx \end{split}$$

$$\begin{split} & \max_{D} V(G,D) \\ & \max_{D} V(G,D) = V(G,D^{*}) \\ & = -2log2 + \int_{x} P_{data}(x)log \frac{P_{data}(x)}{(P_{data}(x) + P_{G}(x))/2} dx \\ & = -2log2 + KL \left(P_{data}||\frac{P_{data} + P_{G}}{2}\right) + KL \left(P_{G}||\frac{P_{data} + P_{G}}{2}\right) \end{split}$$

 $= -2log2 + 2JSD(P_{data}||P_G)$ 

Jensen-Shannon divergence

$$G^* = \arg \min_{G} \max_{D} V(G, D)$$
$$D^* = \arg \max_{D} V(D, G)$$
The maximum objective value is related to JS divergence.



[Goodfellow, et al., NIPS, 2014]

$$G^* = arg \min_{G} \max_{D} V(G, D)$$

$$D^* = arg \max_{D} V(D, G)$$
The maximum objective value is related to JS divergence.
  
• Initialize generator and discriminator
  
• In each training iteration:
$$\underline{Step 1}$$
: Fix generator G, and update discriminator D
$$\underline{Step 2}$$
: Fix discriminator D, and update generator G





- $\theta_G \leftarrow \theta_G \eta \, \partial V(G, D_0^*) / \partial \theta_G$   $\longrightarrow$  Obtain  $G_1$  Decrease JS
- Find  $D_1^*$  maximizing  $V(G_1, D)$  $V(G_1, D_1^*)$  is the JS divergence between  $P_{data}(x)$  and  $P_{G_1}(x)$

• 
$$\theta_G \leftarrow \theta_G - \eta \, \partial V(G, D_1^*) / \partial \theta_G$$
  $\longrightarrow$  Obtain  $G_2$  Decrease JS

divergence(?)

# Algorithm



- Given  $G_0$
- Find  $D_0^*$  maximizing  $V(G_0, D)$

 $V(G_0, D_0^*)$  is the JS divergence between  $P_{data}(x)$  and  $P_{G_0}(x)$ 



### In practice ...

$$V = E_{x \sim P_{data}}[log D(x)] + E_{x \sim P_{G}}[log(1 - D(x))]$$

- Given G, how to compute  $\max_{D} V(G, D)$ 
  - Sample  $\{x^1, x^2, ..., x^m\}$  from  $P_{data}(x)$ , sample  $\{\tilde{x}^1, \tilde{x}^2, ..., \tilde{x}^m\}$  from generator  $P_G(x)$

Maximize 
$$\tilde{V} = \frac{1}{m} \sum_{i=1}^{m} log D(x^i) + \frac{1}{m} \sum_{i=1}^{m} log \left(1 - D(\tilde{x}^i)\right)$$

#### **Binary Classifier**

D is a binary classifier with sigmoid output (can be deep)  $\{x^1, x^2, ..., x^m\}$  from  $P_{data}(x)$  Positive examples  $\{\tilde{x}^1, \tilde{x}^2, ..., \tilde{x}^m\}$  from  $P_G(x)$  Negative examples Minimize Cross-entropy

<b>Algorithm</b> Initialize $\theta_d$ for D and $\theta_g$ for G		
<u>, iigentiii</u>		Can only find $\max V(G, D)$
<ul> <li>In each training iteration:</li> </ul>		lower bound of D
• Sample m examples $\{x^1, x^2,, x^m\}$ from data distribution		
	$P_{data}(x)$	
•	Sample m noise samples $\{z^1, z^2,, z^m\}$ from the prior	
Learning	$P_{prior}(z)$	
D	Obtaining generated data $\{\tilde{x}^1, \tilde{x}^2,, \tilde{x}^m\}, \tilde{x}^i = G(z^i)$	
•	Update discriminator parameters $ heta_d$ to maximize	
<b>Repeat</b> <b>k times</b> • $\tilde{V} = \frac{1}{m} \sum_{i=1}^{m} log D(x^i) + \frac{1}{m} \sum_{i=1}^{m} log \left(1 - \frac{1}{m} \right)\right)\right)\right)\right)\right)$		$-\frac{1}{m}\sum_{i=1}^{m}\log\left(1-D(\tilde{x}^{i})\right)$
	• $\theta_d \leftarrow \theta_d + \eta \nabla \tilde{V}(\theta_d)$	
<ul> <li>Sample another m noise samples {z<sup>1</sup>, z<sup>2</sup>,, z<sup>m</sup>} from t</li> </ul>		
	prior $P_{prior}(z)$	
Learning • Update generator parameters $\theta_{q}$ to minimize		
G Only	• $\tilde{V} = \frac{1}{m} \sum_{l=1}^{m} log D(x^{l}) +$	$-\frac{1}{m}\sum_{i=1}^{m}\log\left(1-D\left(G(z^{i})\right)\right)$
Once	• $\theta_g \leftarrow \theta_g - \eta \nabla \tilde{V}(\theta_g)$	

## Objective Function for Generator in Real Implementation

$$V = E_{x \sim P_{data}} [log D(x)] + E_{x \sim P_{G}} [log (1 - D(x))]$$
  
Slow at the beginning  
Minimax GAN (MMGAN)

$$V = E_{x \sim P_G} \Big[ -log \Big( D(x) \Big) \Big]$$

Real implementation: label x from P<sub>G</sub> as positive

Non-saturating GAN (NSGAN)



# Tips for Improving GAN

Martin Arjovsky, Soumith Chintala, Léon Bottou, Wasserstein GAN, arXiv prepring, 2017 Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, Aaron Courville, "Improved Training of Wasserstein GANs", arXiv prepring, 2017

### JS divergence is not suitable

- In most cases,  $P_G$  and  $P_{data}$  are not overlapped.
- 1. The nature of data

Both  $P_{data}$  and  $P_G$  are low-dim manifold in high-dim space.

The overlap can be ignored.

• 2. Sampling

Even though  $P_{data}$  and  $P_G$  have overlap.

If you do not have enough sampling .....



### What is the problem of JS divergence?

$$P_{G_{0}} \xrightarrow{P_{data}} P_{data} P_{G_{1}} \xrightarrow{P_{data}} P_{data} \cdots P_{G_{100}} P_{data}$$

$$F_{G_{100}} \xrightarrow{P_{data}} P_{data} \xrightarrow{IS(P_{G_{0}}, P_{data})} JS(P_{G_{1}}, P_{data}) \cdots JS(P_{G_{100}}, P_{data})$$

$$= log2 = log2 = 0$$

JS divergence is log2 if two distributions do not overlap.

Intuition: If two distributions do not overlap, binary classifier achieves 100% accuracy

Same objective value is obtained.



Same divergence

# Least Square GAN (LSGAN)



Replace sigmoid with linear (replace classification with regression)



## Wasserstein GAN (WGAN): Earth Mover's Distance

- Considering one distribution P as a pile of earth, and another distribution Q as the target
- The average distance the earth mover has to move the earth.



### WGAN: Earth Mover's Distance



There many possible "moving plans".

Using the "moving plan" with the smallest average distance to define the earth mover's distance.

Source of image: https://vincentherrmann.github.io/blog/wasserstein/

### WGAN: Earth Mover's Distance



There many possible "moving plans".

Using the "moving plan" with the smallest average distance to define the earth mover's distance.

Source of image: https://vincentherrmann.github.io/blog/wasserstein/



A "moving plan" is a matrix The value of the element is the amount of earth from one position to another.

Average distance of a plan  $\gamma$ :

$$B(\gamma) = \sum_{x_p, x_q} \gamma(x_p, x_q) \|x_p - x_q\|$$

Earth Mover's Distance:

 $W(P,Q) = \min_{\gamma \in \Pi} B(\gamma)$ The best plan





[Martin Arjovsky, et al., arXiv, 2017]

### WGAN

Evaluate wasserstein distance between  $P_{data}$  and  $P_{G}$ 

$$V(G,D) = \max_{D \in \underline{1-Lipschitz}} \{ E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[D(x)] \}$$

D has to be smooth enough.

Without the constraint, the training of D will not converge.

Keeping the D smooth forces D(x) become  $\infty$  and  $-\infty$ 



Weight Clipping [Martin Arjovsky, et al., arXiv, 2017] Force the parameters w between c and -c WGAN After parameter update, if w > c, w = c; if w < -c, w = -c Evaluate wasserstein distance between  $P_{data}$  and  $P_{G}$ V(G,D) $= \max_{D \in \underline{1-Lipschitz}} \{ E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[D(x)] \}$ D has to be smooth enough. How to fulfill this constraint? 1–Lipschitz **Lipschitz Function**  $||f(x_1) - f(x_2)|| \le K ||x_1 - x_2||$ Output Input 1-Lipschitz? change change K=1 for "1 - Lipschitz" Do not change fast

### Improved WGAN (WGAN-GP)

$$V(G,D) = \max_{D \in 1-Lipschitz} \{ E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[D(x)] \}$$

A differentiable function is 1-Lipschitz if and only if it has gradients with norm less than or equal to 1 everywhere.

$$D \in 1 - Lipschitz \quad ||\nabla_{x}D(x)|| \leq 1 \text{ for all } x$$

$$V(G, D) \approx \max_{D} \{E_{x \sim P_{data}}[D(x)] - E_{x \sim P_{G}}[D(x)]$$

$$-\lambda \int_{x} \max(0, ||\nabla_{x}D(x)|| - 1)dx\}$$
Prefer  $||\nabla_{x}D(x)|| \leq 1 \text{ for all } x$ 

$$-\lambda E_{x \sim P_{penalty}}[\max(0, ||\nabla_{x}D(x)|| - 1)]\}$$

Prefer  $||\nabla_x D(x)|| \le 1$  for x sampling from  $x \sim P_{penalty}$ 



"Given that enforcing the Lipschitz constraint everywhere is intractable, enforcing it **only along these straight lines** seems sufficient and experimentally results in good performance."

Only give gradient constraint to the region between  $P_{data}$  and  $P_G$  because they influence how  $P_G$  moves to  $P_{data}$
#### Improved WGAN (WGAN-GP)



"Simply penalizing overly large gradients also works in theory, but experimentally we found that this approach converged faster and to better optima."



# Spectrum Norm

Spectral Normalization → Keep gradient norm smaller than 1 everywhere [Miyato, et al., ICLR, 2018]



#### Algorithm of WGAN



[Junbo Zhao, et al., arXiv, 2016]

# Energy-based GAN (EBGAN)

- Using an autoencoder as discriminator D
  - Using the negative reconstruction error of auto-encoder to determine the goodness
  - Benefit: The auto-encoder can be pre-train by real images without generator.



## EBGAN



Auto-encoder based discriminator only gives limited region large value.



Hard to reconstruct, easy to destroy

# Outlook: Loss-sensitive GAN (LSGAN)



## Reference

- Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, Yoshua Bengio, Generative Adversarial Networks, NIPS, 2014
- Sebastian Nowozin, Botond Cseke, Ryota Tomioka, "f-GAN: Training Generative Neural Samplers using Variational Divergence Minimization", NIPS, 2016
- Martin Arjovsky, Soumith Chintala, Léon Bottou, Wasserstein GAN, arXiv, 2017
- Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, Aaron Courville, Improved Training of Wasserstein GANs, NIPS, 2017
- Junbo Zhao, Michael Mathieu, Yann LeCun, Energy-based Generative Adversarial Network, arXiv, 2016
- Mario Lucic, Karol Kurach, Marcin Michalski, Sylvain Gelly, Olivier Bousquet, "Are GANs Created Equal? A Large-Scale Study", arXiv, 2017
- Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, Xi Chen Improved Techniques for Training GANs, NIPS, 2016

## Reference

- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, Sepp Hochreiter, GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium, NIPS, 2017
- Naveen Kodali, Jacob Abernethy, James Hays, Zsolt Kira, "On Convergence and Stability of GANs", arXiv, 2017
- Xiang Wei, Boqing Gong, Zixia Liu, Wei Lu, Liqiang Wang, Improving the Improved Training of Wasserstein GANs: A Consistency Term and Its Dual Effect, ICLR, 2018
- Takeru Miyato, Toshiki Kataoka, Masanori Koyama, Yuichi Yoshida, Spectral Normalization for Generative Adversarial Networks, ICLR, 2018