Introduction to Generative Adversarial Networks (GANs).

HELMHOLTZ RESEARCH FOR GRAND CHALLENGES

CLUSTER OF EXCELLENCE

QUANTUM UNIVERSE

April 29th 2019, Deep Learning and GPU Computing Seminar, DESY Torben Ferber (<u>torben.ferber@desy.de</u>)







Overview

- Introduction to GANs •
- GAN hell: Vanishing gradients and mode collapse •
- WGAN-GP in PyTorch
- GAN metrics: How to measure performance
- GAN hacks •
- Real examples •



Introduction to GANs

- Goal: "synthesise artificial samples, such as images, that are indistinguishable from authentic samples"
- Generative modelling:
 - observe samples to generate more samples •
 - infer density function that describes data
- Several different generative models, GANs are one of them
- GANs introduced to ML community in 2014 by Goodfellow et al. [1], first ideas are much older



[1] GANs (<u>https://arxiv.org/abs/1406.2661</u>)





Edmond de Belamy (\$432,5000)





https://thispersondoesnotexist.com/



StyleGAN (<u>https://arxiv.org/abs/1812.04948</u>)





Generative models: Why?

- Generate simulated environments (e.g. for RL) •
- Infer possible future states •
- Recover missing information •
- Multi-modal outputs (= many possible outputs) •
- Image-to-image translation •
- Cross-domain: Text-to-image •



Generative face completion



(a) Original face (b) Masked input

[1] <u>http://openaccess.thecvf.com/content_cvpr_2017/papers/Li_Generative_Face_Completion_CVPR_2017_paper.pdf</u>

(c) Our result



Introduction to GANs (Torben Ferber) DESY.

Multi-domain image-to-image translation

Input

Blond hair

Aged



Pale skin

Input



Happy Fearful



[1] StarGAN (<u>https://arxiv.org/abs/1711.09020</u>)







GANs in words

Two different models that are adversarial of each other (game theory):

Each model tries to get the highest possible payoff

- Player 1: Generator (G) that tries to capture a training data distribution
- Player 2: Discriminator (D) that tries to identify if data came from the training data or from the generator
- Ideal minimax two player game if G recovers training data and D=0.5 or in other • words: If the two players reach their Nash equilibrium







Introduction to GANs: Nash equilibrium

- game theory:
 - Both players have chosen a strategy •
 - their current strategy

Both strategies and payoffs constitute a Nash equilibrium

Nash equilibrium first formulated by John Nash in the 1950s in the context of

No player can benefit by changing the strategy if the other player sticks to

In other words: A change in strategy will lead to a worse result for this player

THE example: "Prisoners dilemma" (<u>https://en.wikipedia.org/wiki/Prisoner%27s_dilemma</u>)





Introduction to GANs: Schematic

x sampled from data



Introduction to GANs: Schematic



Introduction to GANs: Schematics











http://billhooverart.com



















Introduction to GANs: Training process





pg~p_{data}, D is partially accurate Update D: D is trained until D* = p_{data} / (p_{data} +p_g)

[1] GANs (<u>https://arxiv.org/abs/1406.2661</u>)



Update G:

Gradient of D has guided G to move towards datalike regions Nash equilibrium! D(x) = 0.5

Repeat many times

• • •

GANs: Loss function

binary classifier (D(x)=1 for data and =0 for fake):

$$\max_{D} V(D) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})} [\log I]$$

Objective function for the generator G:

$$\min_{G} V(G) = \mathbb{E}_{z}$$

Minimax game value function:

 $\min \max V(D, G) = \mathbb{E}_{x \sim p_r(x)}[\log D(x)] + \mathbb{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))]$ G

Objective function for the discriminator D is given by the cross entropy of a

 $D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z})))]$

 $\sim_{p_{z}(z)}[\log(1 - D(G(z)))]$



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

- $p_{\text{data}}(\boldsymbol{x}).$
- Update the discriminator by ascending its stochastic gradient:

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

end for

- Update the generator by descending its stochastic gradient:

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

end for

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

• Sample minibatch of m noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$. • Sample minibatch of m examples $\{x^{(1)}, \ldots, x^{(m)}\}$ from data generating distribution

• Sample minibatch of m noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.



The GAN hell

Vanishing gradients

Mode collapse



MinMax or MaxMin: Mode collapse

- $\min_{G} \max_{D} V(G, D)$
 - Good! Optimally the network will converge to pgen = pdata
- $\max_{D} \min_{G} V(G, D)$:
 - than fake
 - distribution

Bad! Generator will converge to a point that is currently considered real rather

If D catches that, G often just moves to the next node instead of the full

Network design does not really specify minmax over maxmin: Mode collapse





DESY. Introduction to GANs (Torben Ferber)

Mode collapse



[1] arxiv:1611.02163



Mode collapse

data

0	0	0	0	0	0	Ø	0	0	0	0	0	0	0	0	0	0	0	0	٥	
1	۱	١	۱	1	1	1	1	ŀ	1)	1)	J	I	١	1	1	1]	
2	Э	г	2	2	2	2	z	2	2	2	2	2	2	J	2	Z	2	2	2	
3	3	3	3	З	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	
Ч	4	Ч	4	۴	4	4	4	4	4	4	4	4	4	4	¥	4	¢	4	ų	
5	5	5	5	5	٢	5	১	5	5	5	5	5	5	5	5	5	5	2	5	
6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	
1	1	7	1	?	1	1	7	1	1	7	7	7)	1	7	Μ	7	7	7	
8	8	8	8	8	8	8	8	8	1	8	8	8	8	8	8	8	ł	8	8	
9	9	9	9	9	9	3	9	9	9	9	9	9	9	9	9	9	9	9	9	

7	1	1	ĺ	ľ	4	1	1	1	1
i	İ	1	i	1	1	1	1	\tilde{L}	ļ
I	1	1	1	1	Ĩ	1	1	ł	1
l	1	2	ł	I	1	1	1		1
1	1	ł	1	1	1	1	Ĵ.	1	1
I	\tilde{L}	1	i	1	1	1	1	1	1
į		1	ł	2	l	ļ	1	3	1
ľ	1	2	ļ	1	7	1	1	ł	1
ļ.	1	Ĵ,	Ŧ	1	}	1	A	1	1
1	1	1	1	Ľ	ſ	ţ	l	1	1

generated



Vanishing gradients: Modified loss function

- Usually D has a much easier start than G, and D can reject all generated samples as fake:
 - $\log(1-D(G(Z)) \rightarrow 0)$
- Slightly modified Minimax game:
 - Instead of training G to minimising log(1-D(G(z)), • train G to maximise log(D(G(z)))
- •

Math behind this: GANs are minimizing the Jenson-Shannon (JS) divergence, which saturates if data distribution **p** and fake distribution **q** do not overlap





Cross-Entropy minus self entropy → Jenson-Shannon (JS) divergence

Maximum likelihood → Kullbach-Leibler (KS) divergence

[1] WGANs 1701.07875







•



[1] <u>https://medium.com/@jonathan_hui/gan-wasserstein-gan-wgan-gp-6a1a2aa1b490</u>

Given the problems (vanishing gradients, mode collapse) of Vanilla GANs, is there a better distance to measure the difference between real and fake?





DESY. Introduction to GANs (Torben Ferber)

Vanishing gradients: Wasserstein distance



+ many more possibilities

Transport plane



EMD is the infimum (greatest lower bound) of all transport planes:

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} \left[\|x - y\| \right]$$

ch-Rubinstein duality \rightarrow

Using the Kantorovic

$$W(\mathbb{P}_r, \mathbb{P}_{\theta}) = \sup_{\|f\|_L \le 1}$$

where f(x) is a 1-Lifschitz function with $|f(x_1)-f(x_2)| \le |x_1 - x_2|$.

- A neural net can learn the function f(x) with two restrictions: •
 - We have to make sure that f(x) is (almost) 1-Lifschitz.
 - We have to train the critic until (almost) convergence.

- $\mathbb{E}_{x \sim \mathbb{P}_r}[f(x)] \mathbb{E}_{x \sim \mathbb{P}_{\theta}}[f(x)]$

[1] WGANs (arXiv:1701.07875v3)







[1] https://de.wikipedia.org/wiki/Lipschitz-Stetigkeit





- Weight clipping (original WGAN):
 - Limit weights to **-c < w < c**.

Gradient penalty (WGAN-GP, Improved WGAN): at most 1 everywhere." \rightarrow penalty term in loss function if norm > 1

$$L = \underbrace{\mathbb{E}_{\tilde{\boldsymbol{x}} \sim \mathbb{P}_{g}} \left[D(\tilde{\boldsymbol{x}}) \right] - \mathbb{E}_{\boldsymbol{x} \sim \mathbb{P}_{r}} \left[D(\boldsymbol{x}) \right]}_{\text{Original critic loss}} \left[D(\boldsymbol{x}) \right] + \lambda \underbrace{\mathbb{E}_{\hat{\boldsymbol{x}} \sim \mathbb{P}_{\hat{\boldsymbol{x}}}}}_{0} \left[(\underline{\boldsymbol{x}}) \right] + \lambda \underbrace{\mathbb{E}_{\hat{\boldsymbol{x}} \sim \mathbb{P}_$$

Clipping value c is hyperparameter: Very hard to optimize

"A differentiable function f is 1-Lipschitz if and only if it has gradients with norm

Computationally expensive Do not use batch norm in critic

 $[\|
abla_{\hat{x}} D(\hat{x}) \|_2 - 1)^2]$

Our gradient penalty



Algo	orithm 1 WGAN, our proposed
the c	default values $\alpha = 0.00005, c = 0$
Req	uire: : α , the learning rate. c_{i}
γ	$n_{\rm critic}$, the number of iterations of
Req	uire: : w_0 , initial critic paramet
1: v	while θ has not converged do
2:	for $t = 0,, n_{\text{critic}} \mathbf{do}$
3:	Sample $\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$ a b
4:	Sample $\{z^{(i)}\}_{i=1}^{m} \sim p(z)$ a
5:	$g_w \leftarrow \nabla_w \left[\frac{1}{m}\sum_{i=1}^{\overline{m}} f_w(x^{(i)})\right]$
6:	$w \leftarrow w + \alpha \cdot \operatorname{RMSProp}(w,$
7:	$w \leftarrow \operatorname{clip}(w, -c, c)$
8:	end for
9:	Sample $\{z^{(i)}\}_{i=1}^{m} \sim p(z)$ a bat
10:	$g_{\theta} \leftarrow -\nabla_{\theta} \frac{1}{m} \sum_{i=1}^{\bar{m}} f_w(g_{\theta}(z^{(i)}))$
11:	$\theta \leftarrow \theta - \alpha \cdot \operatorname{RMSProp}(\theta, g_{\theta})$
12: e	end while

algorithm. All experiments in the paper used $0.01, m = 64, n_{\text{critic}} = 5.$

the clipping parameter. m, the batch size. of the critic per generator iteration. ters. θ_0 , initial generator's parameters.

batch from the real data. batch of prior samples. $(-) - \frac{1}{m} \sum_{i=1}^{m} f_w(g_\theta(z^{(i)}))]$ $(,g_w)$

tch of prior samples.



Algorithm 1 WGAN with gradient penalty. 0.0001, $\beta_1 = 0$, $\beta_2 = 0.9$.
Require: The gradient penalty coefficient λ , n_{critic} , the batch size m , Adam hyperparan
Require: initial critic parameters w_0 , initial g
1: while θ has not converged do
2: for $t = 1,, n_{\text{critic}}$ do
3: for $i = 1,, m$ do
4: Sample real data $x \sim \mathbb{P}_r$, laten
5: $\tilde{\boldsymbol{x}} \leftarrow G_{\theta}(\boldsymbol{z})$
6: $\hat{\boldsymbol{x}} \leftarrow \epsilon \boldsymbol{x} + (1-\epsilon)\tilde{\boldsymbol{x}}$
7: $L^{(i)} \leftarrow D_w(\tilde{\boldsymbol{x}}) - D_w(\boldsymbol{x}) + \lambda(\boldsymbol{x})$
8: end for
9: $w \leftarrow \operatorname{Adam}(\nabla_w \frac{1}{m} \sum_{i=1}^m L^{(i)}, w, c)$
10: end for $m = m = 1$
11: Sample a batch of latent variables $\{z^{(a)}\}$
12: $\theta \leftarrow \operatorname{Adam}(\nabla_{\theta} \frac{1}{m} \sum_{i=1}^{m} -D_{w}(G_{\theta}(\boldsymbol{z})))$
13: end while $(m \sum i = 1)$

[1] WGAN-GPs (arXiv:1704.00028v3)

We use default values of $\lambda = 10$, $n_{\text{critic}} = 5$, $\alpha =$

the number of critic iterations per generator iteration neters α , β_1 , β_2 . generator parameters θ_0 .

Int variable $\boldsymbol{z} \sim p(\boldsymbol{z})$, a random number $\epsilon \sim U[0, 1]$.

$$(\|\nabla_{\hat{\boldsymbol{x}}} D_w(\hat{\boldsymbol{x}})\|_2 - 1)^2$$
$$\alpha, \beta_1, \beta_2)$$

 $\{i^{(i)}\}_{i=1}^{m} \sim p(\boldsymbol{z}).$ $(i), \theta, \alpha, \beta_1, \beta_2)$



DESY. Introduction to GANs (Torben Ferber)

WGANs: Schematics











http://billhooverart.com















Other losses?

GAN Type	
GAN	The o
WGAN	EM di
Improved WGAN	No we
LSGAN	L2 los
RWGAN	Relax
McGAN	Mean
GMMN	Maxir
MMD GAN	Adver
Cramer GAN	Cram
Fisher GAN	Chi-s
EBGAN	Autoe
BEGAN	WGA
MAGAN	Dyna

[1] <u>https://towardsdatascience.com/gan-objective-functions-gans-and-their-variations-ad77340bce3c</u>







(...) improvements can arise from a higher computational budget and tuning more than fundamental algorithmic changes."

tl; dr: Hyperparameters are more important than loss functions. (This is heavily debated and probably problem-dependent.)

"We find that most models can reach similar scores with enough hyperparameter optimization and random restarts.

[1] Google Brain, arXiv:1711.10337v4



















GAN Inference: Only needs the generator





Showcased GAN Output in Papers



and the second se

Actual GAN Output



DESY. Introduction to GANs (Torben Ferber)

ML on HTCondor

1) submit.htc

Universe = vanilla Executable = submit.sh Log = job.log.\$(Cluster)-\$(Process) Error = job.err.\$(Cluster)-\$(Process) Output = job.out.\$(Cluster)-\$(Process) RequestMemory = 4096 Request_GPU = 1 Requirements = (OpSysAndVer == "CentOS7") +RequestRuntime = 18000 notification = Error Queue 1

2) submit.sh

#!/bin/bash
source /etc/profile.d/modules.sh
module load anaconda3

WORK_DIR=`mktemp -d` scp <path_to_my_data_on_pnfs> \$WORK_DIR

python <path_to_my_script>/myscript.py -i traindir \$WORK_DIR

3) myscript.py

import matplotlib.pyplot as plt
plt.switch_backend('agg')

import torch

...



DESY. Introduction to GANs (Torben Ferber)



Or PyTorch






WGAN-GP in PyTorch

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

 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 θ has not converged **do**

for $t = 1, ..., n_{\text{critic}}$ do 2: for i = 1, ..., m do 3: Sample real data $x \sim \mathbb{P}_r$, latent variable $z \sim p(z)$, a random number $\epsilon \sim U[0, 1]$. 4: $\tilde{\boldsymbol{x}} \leftarrow G_{\theta}(\boldsymbol{z})$ 5: $\hat{\boldsymbol{x}} \leftarrow \epsilon \boldsymbol{x} + (1-\epsilon)\tilde{\boldsymbol{x}}$ 6: $L^{(i)} \leftarrow D_w(\tilde{\boldsymbol{x}}) - D_w(\boldsymbol{x}) + \lambda(\boldsymbol{x})$ 7: 8: end for $w \leftarrow \operatorname{Adam}(\nabla_w \frac{1}{m} \sum_{i=1}^m L^{(i)}, w, c)$ 9: end for 10: Sample a batch of latent variables $\{z^{(i)}\}$ 11: $\theta \leftarrow \operatorname{Adam}(\nabla_{\theta} \frac{1}{m} \sum_{i=1}^{m} -D_w(G_{\theta}(\boldsymbol{z})))$ 12: 13: end while

Require: The gradient penalty coefficient λ , the number of critic iterations per generator iteration

$$\begin{array}{l} (\|\nabla_{\hat{x}} D_{w}(\hat{x})\|_{2} - 1)^{2} \\ \alpha, \beta_{1}, \beta_{2}) \\ \\ i) \}_{i=1}^{m} \sim p(\boldsymbol{z}). \\ \beta, \alpha, \beta_{1}, \beta_{2}) \end{array}$$



WGAN-GP in PyTorch: Input \rightarrow Output



Each cell is a 5×5cm Csl(Tl) crystal

1 GeV photons in the Belle II barrel calorimeter (GEANT4)



WGAN-GP in PyTorch: Sample real data

from torch.utils.data import Dataset class ClusterImage(Dataset): 111111 Dataset of cluster images. 111111 def __init__(self, df_x, df_monitor,): self.images = df_x.values.astype(dtype=np.float32) self.monitor = df_monitor.values.astype(dtype=np.float32) self.y = np.ones(len(self.images)) # todo: clip to remove outliers # todo: scale to [-1, 1] using sklearn MinMaxScaler def len (self): return self.y.shape[0] def __getitem__(self, idx): return (self.images[idx], self.monitor[idx], self.y[idx])



DESY. Introduction to GANs (Torben Ferber)

WGAN-GP in PyTorch: Sample real data

```
import torch, uproot, pandas
```

```
# read ROOT input file(s)
tree = uproot.open('/pnfs/desy.de/../myfile.root')['mytree']
df = tree.pandas.df(['var1','var2', 'pixel_*'])
```

```
# create dataset
```

```
# create dataloader
data_loader = torch.utils.data.DataLoader(dataset=dataset,
```







WGAN-GP in PyTorch: Training

```
import time
N_CRITIC_STEPS = 5
for epoch in range(N_EPOCHS):
    # start timer
    t0 = time_time()
    #loop through batches (no need for class labels right now)
    for batch_number, (images, monitor, y) in enumerate(data_loader):
         if critic_step < N_CRITIC_STEPS:</pre>
              # train critic
         else:
              # train generator
    # end timer
    t_epoch = time.time() - t0
    print('time per epoch: %2.2f seconds' % (t_epoch))
```

hyperparameter Important: Train WGANs until (almost) convergence





WGAN-GP in PyTorch: Networks and optimizers







WGAN-GP in PyTorch: Critic and generator training

```
critic.zero_grad()
# clamp weights
# for p in critic.parameters():
                           C = 0.01
      p.data.clamp_(-C, C)
# prediction on data
pred_real = torch.mean(critic(images))
# get latent noise vector
z = to_var(torch_randn(batch_size, N_Z))  N_Z = 10
# generate images
fake_images = generator(z)
# prediction on fake data
pred_fake = torch.mean(critic(fake_images))
# gradient penalty
gradient_penalty = calc_gradient_penalty(critic,
                                         images.data,
                                         fake_images.data,
                                         batch_size,
                                         LAMBDA)
# loss and back propagation
critic_loss = pred_fake - pred_real + gradient_penalty
critic_loss.backward()
critic_optimizer.step()
critic step += 1
```

```
generator.zero_grad()
                  # get latent noise vector
                  z = to_var(torch.randn(batch_size*2, N_Z))
                  # generate images
                   fake_images = generator(z)
                  # prediction on fake data
                   pred_fake = torch.mean(critic(fake_images))
                  # loss and back propagation
                   generator_loss = -pred_fake
                   generator_loss.backward()
                  generator_optimizer.step()
LAMBDA = 10.0
                      def to_var(x):
                          if torch.cuda.is_available():
                              x = x_{\bullet}cuda()
                          return Variable(x)
```





an arrestaring on an - All an - and SPIRISERIUS 10 40 Fr the set of the second second we have a second to the second secon I can't a destate Phasesterstory as appreciated a sur-

:61

a reacta. Carlparaneereda angaina paranaere 200 0.5

0

0





Metrics: Visual

• Works well for real images using real humans to judge image quality





"2-4-legged-Horse-Cow" (Over-generalization)





"Many-eyed-dog". (CNN insensitive)









(Goodfellow 2016)

"The skinned dog." (Orthogonal 2D projection)

[1] Ian Goodfellow 2016





How good is the GAN?

- Losses are difficult:
 - Vanilla GANs play a minimax game where losses are traded
 - WGAN loss is (somewhat) proportional to image quality
- Better:
 - Feature differences like 1D Wasserstein distances (from scipy)
 - Quality measures from image processing:
 - Structural similarity (SSIM) •
 - Kernel principal component analysis (kPCA)
- Classification based metrics (Inception score, ...)



WGAN-GP in PyTorch: Input -> Output

Input



epoch: 00000000







WGAN-GP in PyTorch: Input -> Output

Input



epoch: 0000002







WGAN-GP in PyTorch: Input -> Output

Input





epoch: 00000010





10⁻²

10⁻³

WGAN-GP in PyTorch: Input -> Output

Input











WGAN-GP in PyTorch: Input -> Output

Input



epoch: 00002000







How good is the GAN?





















Precision, recall, and F1

- Quality indicator for (multi) mode discriminative models.
 - **Precision**: $P = T_p / (T_p + F_p)$ •
 - **Recall**: $R = T_p / (T_p + F_n)$ •
 - **F1**: $F1 = 2 \times (P \times R) / (P + R)$, Bad is 0, Good is 1 •
- Be aware that relative importance of precision and recall is very problem dependent:
 - Assume you want to detect nuclear missile attacks and have one false negative...







PRD Curves (Precision and recall for distributions)



[1] "Assessing Generative Models via Precision and Recall" <u>https://arxiv.org/abs/1806.00035</u> [2] <u>https://github.com/msmsajjadi/precision-recall-distributions</u>





Inception score (IS)

- ImageNet)
- A well performing GAN generator will produce: •
 - •
 - •
- IS = Exponential of KL distance between class distribution and marginal distribution
 - High score is better (unbound)

Use a powerful and well-trained discriminator (originally Inception-v3 on

Clear images for every mode (low entropy for conditional label distribution)

All modes, correct fraction of modes (high entropy for marginal distribution)

[1] arxiv:1606.03498







Inception score (IS)

Similar labels sum to give focussed distribution





[1] <u>https://cdn-images-1.medium.com/max/2400/1*23gj_d3dxfm5FoKae_pc5Q.png</u>

Different labels sum to give uniform distribution







Inception score (IS): Potential problems

- IS depends on (powerful) classifier and their training data
- necessarily what one desires.
- Generator overtraining will not be penalized •
- Low variance of images per mode will not be penalized

→ Fréchet Inception Distance (FID) is generally superior (but more • complicated): Small FID is better

Classifiers often rely on CNNs which generally rely on local textures. This is not





There is not one single metric to compare GANs.









One sided-label smoothing (for Vanilla GANs)

- Goal: Avoid overconfidence of Discriminator
- Almost trivial idea, but very powerful.
- Smooth labels for real images $(1 \rightarrow 0.9)$: Reduce confidence for real images to • reduce steep gradients.
- Never (!) smooth labels for fake images ($0 \rightarrow 0.1$): This would encourage the generator to continue producing wrong images.
- A variant is to smooth with a random value instead of a fixed 0.1.

[1] <u>https://arxiv.org/abs/1701.00160</u>





Instance and label noise

- Goal: Create overlapping distributions •
- Add noise to generated and real images before passing them to generator (instance noise)
- and/or flip labels of real images (label noise)
- Can stabilise early training but will • reduce image quality

def add_instance_noise(

noise = Variable(std return images + nois



[1] <u>https://www.inference.vc/</u>





Minibatch discrimination

- Goal: Avoid mode collapse •
- Measure similarity of all images in a batch. If similarity increases: Mode collapse
- Add similarity measure(s) as input to (first) fully connected layer
- Can be computationally intensive if • batch size is large
- Poor man's version: Just detect if images are the same by checking one statistics, e.g. the mean

```
class critic(nn.Module):
# ...
    def main(self, x):
        out = torch.cat((x,
                 self.mbd(x).cuda()), dim=1)
        for layer in self.layer_module:
            out = layer(out)
        return out.view(out.size(0), -1)
```









Batch normalization

- Goal: Stable training (used in DCGAN and generally in every CNN) •
- Problem: Introduces intra-batch correlations



[1] GANs (<u>https://arxiv.org/abs/1406.2661</u>)

Solution: Virtual Batch Norm using a reference batch chosen before training



Experience replay

- Keep a history of generated images
- Feed a (small) fraction of old images to discriminator to avoid overtraining of the discriminator

GUIDE TO MAKING PEOPLE FFELOLD IF THEY'RE [AGE], YOU SAY: "DID YOU KNOW [THING] HAS BEEN AROUND FOR A MAJORITY OF YOUR LIFE?" THING AGE - GRAND THEFT AUTO 🛽 16 - RICKROLLING AQUA TEEN HUNGER FORCE COLON MOVIE FILM FOR THEATRES - THE NINTENDO WII - TWITTER 20 — - THE XBOX 360, XKCD - CHUCK NORRIS FACTS 22 — - OPPORTUNITY'S MARS EXPLORATION - FACEBOOK 24 — - GMAIL, PIRATES OF THE CARIBBEAN — IN DACLUB 26 -— FIREFLY 27 — - THE WAR IN AFGHANISTAN 29 ----- THE iPOD 30 ----- SHREK, WIKIPEDIA 3) — 32 ----- THE SIMS 33 - AUTOTUNED HIT SONGS -THE STAR WARS PREQUELS 34 — — THE MATRIX 35 ----- POKÉMON RED& BLUE 36 -- NETFLIX, HARRY POTTER, GOOGLE 37 — 38 ----- DEEP BLUE'S VICTORY 39 ----- TUPAC'S DEATH - THE LAST CALVIN AND HOBBES STRIP 40 -TOY STORY ----- [DON'T WORRY, THEY'VE GOT THIS COVERED] >41

THE NOVEMBER 2016

[1] <u>https://xkcd.com/1757/</u>





Historical averaging

- average)
- from past average:

Keep history of model parameters θ for the last t models (or keep the running

Add L2 penalty to loss function if current model parameters deviate too much









Class conditioning

- •
- bad perceived image quality, ...
- Possible solution is called class conditioning: •
 - Class-conditional (add class to latent space of generator)
 - Class discrimination (add class discrimination to discriminator) •
 - Combinations of those •

Standard (W)GANs have no additional information about the modes in data

Several problems in GANs are related to that: Slow training, mode collapse,



Class conditioning



[1] <u>https://github.com/znxlwm/pytorch-generative-model-collections</u>













[1] "CaloGAN: Simulating 3D High Energy Particle Showers in Multi-Layer Electromagnetic Calorimeters with Generative Adversarial Networks" arXiv:1712.10321 (LAGAN: arXiv:1701.05927)





DESY. Introduction to GANs (Torben Ferber)

LHCbGAN



[2] "Generative Models for Fast Calorimeter Simulation: the LHCb case" arXiv:1812.01319v2



Summary

- GANs are cool, fun, and sometimes almost magical •
- GANs are difficult
- GANs are a young field in ML with ongoing basic research •
- - for physics
 - Get in touch to learn from each other •

Numerous applications in physics (simulation, calibration, understanding, ...)

Many things that are irrelevant when creating human faces are very relevant





Summary

- GANs are cool, fun, and sometimes almost magical •
- GANs are difficult
- GANs are a young field in ML with ongoing basic research •
- - for physics
 - Get in touch to learn from each other •

Numerous applications in physics (simulation, calibration, understanding, ...)

Many things that are irrelevant when creating human faces are very relevant




Contact

DESY.

Deutsches Elektronen Synchrotron <u>www.desy.de</u> Torben Ferber <u>torben.ferber@desy.de</u> ORCID: 0000-0002-6849-0427