

Generative Models: Part II

Advanced GAN Techniques & Application in Particle Physics

Wasserstein GANs

Jonas Glombitza, Martin Erdmann

RWTH Aachen

2nd Terascale Machine Learning School |



Outline

I. Variational Autoencoders

II. Introduction to Generative Adversarial Networks (GANs)

> Adversarial frameworks

III. Tutorial: Implementation of GANs

BREAK

IV.Latest developments & advanced techniques

- Wasserstein GANs
- V.Application in physics research
 - Simulation acceleration
 - Style transfer (domain adaption)

VI.Tutorial: Implementation of Wasserstein GANs





Feel free to ask questions during the seminar! Just "raise" your hand...

Recap – Generative Models



Variational Autoencoder

- Hidden representation follows low dimensional arbitrary prior distribution
- Trained decoder part can be used as generator to produce new samples





Generative Adversarial Networks

- Hand-coded loss is replaced by discriminator (tries to discriminate between fake samples and real samples)
- Adversarial training: generator and discriminator trained "against" each other
- Generator tries to fool discriminator

Tutorial on Generative Models

Recap - Manifold Hypothesis



Idea: Manifolds of meaningful pictures are highly concentrated with very little volume and embedded in a very high dimensional space

- Generation of images is a very challenging task
- Correlations / probability dimension are high dimensional
- **Example:** Try to generate images randomly:



Goal







Sample 1

Sample 100,000

You will even never reach this "neighborhood sample"

"To deal with a 14-dimensional space, visualize a 3-D space and say 'fourteen' to yourself very loudly. Everyone does it." - G. Hinton

Tutorial on Generative Models

Results





Tutorial on Generative Models

5

Results Fashion MNIST



- Challenging training!
- "Weak convergence" after 3000 iterations
- Complex models show very instable training

Tutorial on Generative Models

6







Interpreting the Adversarial Loss

- GANs are hard to train \rightarrow Nash equilibrium
 - generator \longleftrightarrow discriminator

- Loss is hard to interpret (depends on discriminator)
 - no correlation with image quality

- Strong discriminator \rightarrow vanishing gradients
- Best: generator and discriminator on same scale
 - Inexact noisy training \rightarrow Rarely converging framework

Tutorial on Generative Models

Glombitza | RWTH Aachen | 03/05/21 | Part 10: Advanced GAN Techniques & Application in Particle Physics





Mode Collapsing - Helvetica Scenario

Problem: GANs often suffer from mode collapsing

- Many $\mathbf{z} \sim p(\mathbf{z})$ collapse towards restricted space in P_r
 - Generator produce samples of a limited phase space
 - Example: generate only digits 1 and 8
- Discriminator feedback is insensitive to complete phase-space
 - Will focus on point(s) of phase space the generator do not cover
- Discriminator will push generator to this mode \rightarrow cycling behavior
- Need different (softer) metric to address these issues!







GAN Objective

RWTHAACHEN UNIVERSITY

• By fooling the discriminator the generator minimize distribution differences $\rightarrow P_{\theta} \approx P_r$



1

• GAN training similar to minimizing Jensen-Shannon divergence (assume optimal discriminator)

$$\mathcal{D}_{JS}(P_r||P_\theta) = \mathcal{D}_{KL}(P_r||P_m) + \mathcal{D}_{KL}(P_\theta||P_m) \qquad P_m = \frac{1}{2}(P_r + P_\theta)$$

- Symmetrized and smoothed version of the Kullback-Leibler divergence
- \boldsymbol{x} Fails to provide a meaningful value when two distributions are disjoint
 - In very high dimensional manifolds the distributions between generated and real samples are disjoint

Tutorial on Generative Models

9

Wasserstein Distance



- Also known as Earth Mover's distance (EMD) Ensures smallest cost $\mathcal{D}_W(P_r||P_\theta) = \inf_{\substack{\gamma \in \Pi(P_r, P_\theta)}} \mathbb{E}_{(x,y) \sim \gamma}[||x - y||]$ Transportation plans
- Describes minimal cost to move distribution P_{θ} on P_r and vice versa
 - Cost: mass * distance



Tutorial on Generative Models

Distribution Similarity - Metrics

- Kullback-Leibler divergence
 - × Not finite, not symmetric

$$\mathcal{D}_{KL}(P_r||P_{\theta}) = \mathbb{E}_{\mathbf{x} \sim P_r} log\left(\frac{P_r}{P_{\theta}}\right)$$

- Jensen-Shannon divergence $\mathcal{D}_{JS}(P_r||P_{\theta}) = \mathcal{D}_{KL}(P_r||P_m) + \mathcal{D}_{KL}(P_{\theta}||P_m) \qquad P_r$ • Symmetric
 - $P_m = \frac{1}{2}(P_r + P_\theta)$

For disjoint distributions: $\mathcal{D}_{KL}(P_{\theta}||P_r) = \infty$

 $\mathcal{D}_{JS}(P_r||P_\theta) = \log(2)$

 $\mathcal{D}_{KL}(P_r||P_\theta) = \infty$

Symmetric

Wasserstein distance

Meaningful distance measure for disjoint distributions

In GAN training we are dealing with disjoint distributions!

Kantorovich-Rubinstein Duality

• Wasserstein distance (formal definition very intractable) can be expressed

$$\mathcal{D}_W(P_r||P_\theta) = \sup_{f \in Lip_1} \mathbb{E}_{x \sim P_r}[f(x)] - \mathbb{E}_{\tilde{x} \sim P_\theta}[f(\tilde{x})]$$

- supremum = least upper bound
- f =Set of 1-Lipschitz functions
- $\mathbb{E}_{x \sim P_r}[f(x)]$ Expectation value when applying set of 1-Lipschitz functions on samples from real samples





Slope everywhere less equal 1!

> Approximate $f \approx f_w$ with neural network



The WGAN Concept



- Neural network carries the Lipschitz continuity constraint
- Critic network estimate Wasserstein distance between generate and real samples

To paint more realistic images: Just change your brush!



Tutorial on Generative Models

13

Gradient Penalty

- Implement Lipschitz constraint
- > Build up space for meaningful discriminator feedback
- Without Lipschitz constrain
 - Critic will not converge → No Wasserstein!

Extend objective with additional term:

- Penalize gradients being different from 1 $\mathcal{L}_{GP} = \lambda \mathbb{E}_{\hat{u} \sim P_{\hat{u}}} [(||\nabla_{\hat{u}} f_w(\hat{u})||_2 - 1)^2]$ • hyperparameter
- Sample gradients along line between event mixture \hat{u}

$$\hat{u} = \epsilon x + (1 - \epsilon)\tilde{x} \qquad 0 \le \epsilon \le 1$$













Critic

- Critic approximates Wasserstein distance
 - Carries Lipschitz constraint
 - Ensures meaningful and stable gradients
- No explicit formulation of loss function
 - Approximate loss function itself
 - Maximize difference $|f_w(\tilde{x}) f_w(x)|$
- Critic should be always trained to convergence
 - Usually ~ 5 10 iterations

$$\mathcal{D}_W(P_r||P_\theta) = \sup_{f \in Lip_1} \mathbb{E}_{x \sim P_r}[f_w(x)] - \mathbb{E}_{\tilde{x} \sim P_\theta}[f_w(\tilde{x})]$$



Advantages WGAN vs. GAN



- Train critic to convergence ensure quality gradients
- Insensitive to mode collapsing

16

- Meaningful metric / objective \rightarrow allow for easy hyperparameter search
- Convergence correlates with generation quality
- Change from Jensen Shannon divergence to Wasserstein-1
- We get feedback from an art expert!



Results

III. Physikalisches III. Physikalisches UNIVERSITY

- WGAN generates images with much better quality
- Critic loss converges
- Loss correlates with images quality



Wasserstein GANs

- Allow stable training of GANs
 - Train critic to convergence
 - Precise feedback for generator
- Prevent mode collapsing
- Provide meaningful loss

Tutorial on GANs, IML Workshop 2019, CERN Glombitza | RWTH Aachen | 03/05/21



Application in Particle Physics



- Detector simulation are very time consuming
- Replace simulation programs like Geant4 with generative model
 - Reach speed-up of factor 10³ 10⁵
- Add constrainer networks to condition the generation process
 - Generator needs dependence (energy, particle type...)
 - Samples must comply with physics laws





Paganini, Oliviera, Nachman - https://arxiv.org/abs/1712.10321

Erdmann, Geiger, Glombitza, Schmidt - https://arxiv.org/abs/1802.03325

Generation of Calorimeter Images



- Quality of images is crosschecked using physics observables
- Challenges: Sparsity, logarithmic intensity distribution



Glombitza | RWTH Aachen | 03/05/21 | 2nd Terascale Machine Learning School



GAN

Laver



GAN applications for fast calorimeter simulation in other experiments ²²



Thorben Quast | Auger ML Days, 04 Nov 2020



visible cell energy [MeV]

number of hits

Simulation Refinement

- Simulation data mismatches
- Predictive models can be sensitive to artifacts / mismatches existing in simulation ≻

- Can lead to reconstruction errors.
- Use adversarial networks with refinement constraint
- Train refiner to refine simulations



Tutorial on Generative Models







Simulation Refinement



• Use auto encoder set up to mitigate data / simulation differences





- Simulation and data share encoder but different decoder (similar representation)
- After training: refine simulation with decoder trained to reconstruct data



 $Chintan\ {\tt Trivedi:}\ {\tt Using}\ {\tt Deep}\ {\tt Learning}\ to\ {\tt improve}\ {\tt FIFA}\ {\tt 18}\ {\tt graphics}\ {\tt - \ {\tt Towards}\ {\tt Data}\ {\tt Science}$



Spectral Normalization for GANs



- Gradient penalty / regularization is most important for training GANs!
- WGAN-GP is state of the art
- Adapt Lipschitz constraint using Gradient "normalization" (penalty)
 - Also standard (NS-GAN) with gradient penalty performs well!

- Adapt Lipschitz constraint in the weights using the *spectral norm*
- Critic in WGAN-GP needs many iterations \rightarrow slow training
- Spectral norm can be fast approximated using power iteration method
- Increased stability (high learning rates, high momentum rates)

Tutorial on Generative Models

23



Spectral Norm

24

• Spectral norm: "natürliche Matrixnorm"

$$||\mathbf{W}||_{2} = \max_{x \neq 0} \frac{||\mathbf{W}x||_{2}}{||x||_{2}} = \max_{||x||_{2}=1} ||\mathbf{W}x||_{2} = \sqrt{\lambda_{max}}$$

- Maximum stretch factor of unit vector after multiplication with matrix
- $\lambda_1 = \lambda_{max}$ = highest singular value ("Singulärwert") of the matrix



Spectral Normalization for GANs



 $\mathbf{I}\mathbf{I}\mathbf{I}$

- D(x) = discriminator
- Adapt WGAN-GP constraint (gradient wrt. \boldsymbol{x} real and fake samples)
 - Use **spectral normalization** in each layer!
- Basic idea:

25

$$||D(x)||_{\text{Lip}} = \sup_{x} \sigma(\nabla_x D(x)) = \sup_{x} \sigma(\nabla_x Wx) = \sigma(W) \longrightarrow W_{\text{norm}} = \frac{w}{\sigma(W)}$$

Cover Lipschitz constraint by normalizing the weights

• Gradient update:

• Gradient penalizes updates in direction of highest singular value (in each layer)

Tutorial on Generative Models



SNGAN

discriminator

- discriminator weights cover Lipschitz constraint due to spectral norm
- "regularize" gradients \rightarrow mode collapsing unlikely
- discriminator loss still meaningless \rightarrow no critic / distance measure

generator

- Also spectral normalization in generator improves stability
- > enforce harmless mapping
- Framework trained with 1:1 discriminator / generator update ratio

Game: Which face is real?



http://www.whichfaceisreal.com/index.php





Changing the Generator enerator 000

Progressive Growing

RWTHAACHEN UNIVERSITY

- Separation of training process in several steps
- Increase image resolution stepwise
- Beginning: (low resolution) data set has only few modes
 - > small differences to be learned
- low resolution
 - learn large scale structure
- High resolution
 - learn fine details
- Speed up
 - Most iteration in the beginning
- Tutorial on Generative Models
- ²⁹ Glombitza | RWTH Aachen | 03/05/21 | 2nd Terascale Machine Learning School



Progressive Training





Training on 8 Tesla V100!

Tutorial on Generative Models

30



StyleGAN

- Image contains of several style levels
- > Change structure of generator
 - disentangle styles in architecture
 - Coarse (pose, face shape)
 - Medium (facial features, eyes closed)
 - Fine (finer hair details, exes)
 - Learning of high-level attributes
- Add additional noise
- Use medium representation of latent space
 - Use mapping network





High resolution: 1024 x 1024 pixels



StyleGAN

- Mapping network learns "medium" latent space
- Generator input:
 - Styles control Adaptive instance normalization

$$AdaIN(\mathbf{x}_{i}, \mathbf{y}) = \mathbf{y}_{s,i} \frac{\mathbf{x}_{i} - \mu(\mathbf{x}_{i})}{\sigma(\mathbf{x}_{i})} + \mathbf{y}_{b,i}$$

- Styles **re-scale features maps** in generator at each representation level
- Additional noise at each level of representation
- Sample from restricted (truncated) phase space
 - Loss off variation \rightarrow better quality





Video: A Style based Generator



https://www.youtube.com/watch?time_continue=370&v=kSLJriaOumA

References & Further Reading



- Goodfellow et al.: Generative Adversarial Networks https://arxiv.org/abs/1406.2661
- Arjovsky, Chintala, Bottou: Wasserstein GAN https://arxiv.org/abs/1701.07875
- Gulrajani et al.: Improved Training of Wasserstein GANs https://arxiv.org/abs/1704.00028
- Paganini, Oliveira, Nachman: CaloGAN https://arxiv.org/abs/1712.10321
- Erdmann, Geiger, Glombitza, Schmidt https://arxiv.org/abs/1802.03325
- Erdmann, Glombitza, Quast: Calorimeter WGAN T. Comput Softw Big Sci (2019) 3: 4
- C. Trivedi: Using Deep Learning to improve FIFA 18 graphics Towards Data Science
- Emanuele Sansone: https://github.com/emsansone/GAN
- Miyato et al.: SNGAN- https://arxiv.org/abs/1802.05957
- T. Karras, S. Laine, T. Aila: A Style-Based Generator https://arxiv.org/abs/1812.04948

Tutorial on Generative Models



Summary

GANs

• Delicate, hard to train (mode collapsing, vanishing gradients, meaningless loss)

Advanced techniques – WGAN, ProGAN, SNGAN

- Stabilize training process using smooth metric \rightarrow meaningful distance measure
- Use regularization: penalize gradient, enforce spectral norm of the weights
- Imply prior on generator architecture (progressive growing, hierarchy control)

Generative Models in Physics Research

- Speed up simulations by factor $10^3 10^5$
- Reduce data / simulation mismatches



Tutorial on Generative Models



Tutorial

Open jupyter notebooks in google colab: You can find the repository at: https://github.com/Napoleongurke/tutorial_generative_models

Open PART I: Vanilla_GAN.ipynb



Tutorial on Generative ModelsGlombitza | RWTH Aachen | 03/05/21





Stay tuned...



Yang, Chou, Yang - https://arxiv.org/abs/1703.10847





Zhu, Park, Isola, Efros - https://arxiv.org/abs/1703.10593



Figure 2: From left to right: bicubic interpolation, deep residual network optimized for MSE, deep residual generative adversarial network optimized for a loss more sensitive to human perception, original HR image. Corresponding PSNR and SSIM are shown in brackets. [$4 \times$ upscaling]

Ledig et. al. - https://arxiv.org/abs/1609.04802



Isola, Zhu, Zhou, Efros - https://arxiv.org/pdf/1611.07004.pdf

... there is much more going on!

37 Tutorial on Generative Models

The WGAN-GP Algorithm



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

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 θ has not converged do

for $t = 1, ..., n_{\text{critic}}$ do 2: 3:

for
$$i=1,...,m$$
 do

Sample real data $\boldsymbol{x} \sim \mathbb{P}_r$, latent variable $\boldsymbol{z} \sim p(\boldsymbol{z})$, a random number $\epsilon \sim U[0, 1]$.

5:
$$\tilde{\boldsymbol{x}} \leftarrow G_{\theta}(\boldsymbol{z})$$

6:
$$\hat{x} \leftarrow \epsilon x + (1 - \epsilon)\tilde{x}$$

7: $L^{(i)} \leftarrow D_w(\tilde{x}) - D_s$

$$L^{(i)} \leftarrow D_w(\tilde{\boldsymbol{x}}) - D_w(\boldsymbol{x}) + \lambda(\|\nabla_{\hat{\boldsymbol{x}}} D_w(\hat{\boldsymbol{x}})\|_2 - 1)^2$$

end for 8:

9:
$$w \leftarrow \operatorname{Adam}(\nabla_w \frac{1}{m} \sum_{i=1}^m L^{(i)}, w, \alpha, \beta_1, \beta_2)$$

11: Sample a batch of latent variables
$$\{z^{(i)}\}_{i=1}^m \sim p(z)$$
.

12:
$$\theta \leftarrow \operatorname{Adam}(\nabla_{\theta} \frac{1}{m} \sum_{i=1}^{m} -D_w(G_{\theta}(\boldsymbol{z})), \theta, \alpha, \beta_1, \beta_2)$$

13: end while

4:

Gulrajani et al.: Improved Training of Wasserstein GANs - https://arxiv.org/abs/1704.00028

Tutorial on Generative Models

Generator

- Use critic feedback to increase generation quality
 - Minimize $\mathcal{D}_W(P_r, P_\theta)$ using gradient descent

$$\nabla_{\theta} \mathcal{D}_W(P_r, P_{\theta}) = -\mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\nabla_{\theta} f(G_{\theta}(\mathbf{z}))]$$

No vanishing gradients for the generator





Tutorial on Generative Models

Non Saturation GAN (NS-GAN)



- Use **label switching** to avoid vanishing gradients in discriminator
- Standard loss: *minimize*

40

 $Loss = \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [log(1 - D(G_{\theta}(\mathbf{z})))]$

- But gradients vanish for $D(G_{\theta}(\mathbf{z}))
 ightarrow 0$ (good discriminator)
- Replace loss and minimize instead

$$Loss = -\mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})}[log(D(G_\theta(\mathbf{z})))]$$

No vanishing gradient but very instable update

• But new loss has strange update behavior:

$$\mathbb{E}_{z \sim p(z)} \left[-\nabla_{\theta} \log D^*(g_{\theta}(z)) |_{\theta = \theta_0} \right] = \nabla_{\theta} \left[KL(\mathbb{P}_{g_{\theta}} \| \mathbb{P}_r) - 2JSD(\mathbb{P}_{g_{\theta}} \| \mathbb{P}_r) \right] |_{\theta = \theta_0}$$

Distribution Similarity - Metrics





> Only \mathcal{D}_W provides meaningful distance measure even for disjoint distributions!

Tutorial on Generative Models

41

- Heavily constraints the discriminator
- Gradient Penalty allows for a much more complex approximation

Tutorial on Generative Models

42

Glombitza | RWTH Aachen | 03/05/21 | 2nd Terascale Machine Learning School

Weight Clipping vs. WGAN-GP

- Weight Clipping:
 - Constraints the weights to lie on a compact space
 - Clip weights after each gradient update eg. to [-0,001; 0,001]



-0.02

-0.01

Weight clipping

0.00

Weights

0.01



0.02 - 0.50

-0.25





0.25

0.50





Gradient penalty

0.00

Weights

Progressive Growing

- New layer act as "residual block"
 - In generator & discriminator
- During transition:
 - New block is slowly faded in
 - α increases linear from 0 to 1
- Training samples:
 - downscaled and interpolated during transition between resolutions

D

Tutorial on Generative Models

43 Glombitza | RWTH Aachen | 03/05/21 | 2nd Terascale Machine Learning School



fromRGB

16x16

(a)





1-α 🗸 🕇 α

fromRGB

32x32

0.5x

 $1 - \alpha + \alpha$

16x16

(b)

0.5×

fromRGB