Generative Models II

Aaron Courville CIFAR Fellow, Université de Montréal

CIFAR-CRM Deep Learning Summer School

Université de Montréal, June 29th, 2017

Institut des algorithmes d'apprentissage de Montréal







- Genera distribu
- Density

• Sample

• Density estimation



Sample lacksquare

















Generative models II: Outline

- Autoregressive models
 - PixelCNN
- Latent variable models
 - Variational Autoencoders
 - VAE
 - Improved strategies for inference
 - Generative Adversarial Networks
 - GAN
 - Wasserstein GAN
 - ALI





Autoregressive generative models

- Autoregressive generative models are well known for sequence data \bullet (language modeling, time series, etc.)
- Less obviously applicable to arbitrary (non-sequential) observations
- Some history: logistic regression for the conditionals (Frey et al., 1996) neural networks for the conditionals (Bengio and Bengio, 2000) idem, with new weight sharing (NADE) (Larochelle and Murray, Gregor and Lecun, 2011) Deep NADE, PixelRNN, PixelCNN, WaveNet, Video Pixel Network, etc.







Autoregressive generative models

- Choose an ordering of the dimensions in x.
- Define the conditionals in the product rule expression of p(x).

 $p(\mathbf{x}) = \mathbf{T}$ k =

- Properties
 - Pros: p(x) is tractable, so easy to train, easy to sample (though slower)
 - Cons: doesn't have a natural latent representation





$$\int_{=1}^{0} p(x_k | \mathbf{x}_{< k})$$



Idea: use masked convolutions to e



Oord, Aaron van den, Nal Kalchbrenner,









Oord, Aaron van den, Nal Kalchbrenner, and Koray Kavukcuoglu. "Pixel recurrent neural networks." arXiv preprint arXiv:1601.06759 (2016).







How can convolutions make this raster scan faster?



Training can be parallelized, though generation is still a sequential operation over pixels

van den Oord, Aaron, et al. "Conditional image generation with PixelCNN decoders." Advances in Neural Information Processing Systems. 2016. 8



Use a stack of masked convolutions











Oord, Aaron van den, Nal Kalchbrenner, and Koray Kavukcuoglu. "Pixel recurrent neural networks." arXiv preprint arXiv:1601.06759 (2016). 9



composing multiple layers increases the context size

only depends on pixel above and to the left

masked convolution





Improving PixelCNN

There is a problem with this form of masked convolution.

1	1	1	1	1
1	1	1	1	1
1	1	0	0	0
0	0	0	0	0
0	0	0	0	0

van den Oord, Aaron, et al. "Conditional image generation with PixelCNN decoders." Advances in Neural Information Processing Systems. 2016. 10







Stacking layers of masked convolution creates a blindspot





Improving PixelCNN I



Stacking layers of masked convolution creates a blindspot

van den Oord, Aaron, et al. "Conditional image generation with PixelCNN decoders." Advances in Neural Information Processing Systems. 2016. 11





Solution: use two stacks of convolution, a vertical stack and a horizontal stack





Improving PixelCNNI

This information flow (between vertical and horizontal stacks) preserves the correct pixel dependencies

Split feature maps

van den Oord, Aaron, et al. "Conditional image generation with PixelCNN decoders." NIPS 20162



Use more expressive nonlinearity: $\mathbf{h}_{k+1} = \tanh(W_{k,f} * \mathbf{h}_k) \odot \sigma(W_{k,q} * \mathbf{h}_k)$









PixelCNN: Experimental Results

Topics: CIFAR-10

Samples from a class-conditional PixelCNN

















Conditional Image Generation with PixelCNN Decoders van den Oord, Kalchbrenner, Vinyals, Espeholt, Graves, Kavukcuoglu, NIPS 2016































































































Conditional Image Generation with PixelCNN Decoders van den Oord, Kalchbrenner, Vinyals, Espeholt, Graves, Kavukcuoglu, NIPS 2016



Sorrel horse





14

PixelCNN: Experimental Results

Topics: CIFAR-10

Samples from a class-conditional PixelCNN





Conditional Image Generation with PixelCNN Decoders van den Oord, Kalchbrenner, Vinyals, Espeholt, Graves, Kavukcuoglu, NIPS 2016





PixelCNN: Experimental Results

Topics: CIFAR-10

Performance measured in bits/dim

Model

Uniform Distribution: [30] Multivariate Gaussian: [30] NICE: [4] Deep Diffusion: [24] DRAW: [9] Deep GMMs: [31, 29] Conv DRAW: [8] RIDE: [26, 30] PixelCNN: [30] PixelRNN: [30]

Gated PixelCNN:

NL]



Conditional Image Generation with PixelCNN Decoders van den Oord, Kalchbrenner, Vinyals, Espeholt, Graves, Kavukcuoglu, NIPS 2016

L Test (Train)
8.00
4.70
4.48
4.20
4.13
4.00
3.58 (3.57)
3.47
3.14 (3.08)
3.00 (2.93)
3.03 (2.90)





Parallel Multiscale Autoregressive Density Estimation

Scott Reed, Aaron vanden Oord, Nal Kalchbrenner, Sergio Go'mez Colmenarejo, Ziyu Wang, Dan Belov, Nando de Freitas (2017)

Can we speed up the generation time of PixelCNN?

• Yes, via multiscale generation:





Parallel Multiscale Autoregressive Density Estimation

Scott Reed, Aaron vanden Oord, Nal Kalchbrenner, Sergio Go'mez Colmenarejo, Ziyu Wang, Dan Belov, Nando de Freitas (2017)

Can we speed up the generation time of PixelCNN?

- Yes, via multiscale generation.
- Also seems to help to provide better global structure

"A yellow bird with a black head, orange eyes and an orange bill."





- The Variational Autoencoder model:
 - Representations (ICLR) 2014.
 - latent Gaussian models. ICML 2014.



Image from: Ward, A. D., Hamarneh, G.: 3D Surface Parameterization Using Manifold Learning for Medial Shape Representation, Conference on Image Processing, Proc. of SPIE Medical Imaging, 2007 19





Kingma and Welling, Auto-Encoding Variational Bayes, International Conference on Learning

Rezende, Mohamed and Wierstra, Stochastic back-propagation and variational inference in deep

Frey Face dataset:

 Z_2



Expression

Pose

 z_1



0002 660000 9999777111 7777771111111

 z_2 '

MNIST:

 z_1



latent variable model: learn a mapping from some latent variable z to a complicated distribution on x.

$$p(x) = \int p(x, z) \, dz$$

Can we learn to decouple the true explanatory factors underlying the data distribution? E.g. separate identity and expression in face images



Image from: Ward, A. D., Hamarneh, G.: 3D Surface Parameterization Using Manifold Learning for Medial Shape Representation, Conference on Image Processing, Proc. of SPIE Medical Imaging, 2007 21



where
$$p(\boldsymbol{x}, \boldsymbol{z}) = p(\boldsymbol{x} \mid \boldsymbol{z})p(\boldsymbol{z})$$

p(z) = something simple $p(x \mid z) = g(z)$

latent variable model: learn a mapping from some latent variable z to a complicated distribution on x.

$$p(x) = \int p(x, z) \, dz$$

• Can we learn to decouple the true explanatory factors underlying the data distribution? E.g. separate identity and expression in face images



Image from: Ward, A. D., Hamarneh, G.: 3D Surface Parameterization Using Manifold Learning for Medial Shape Representation, Conference on Image Processing, Proc. of SPIE Medical Imaging, 2007 22



where
$$p(\boldsymbol{x}, \boldsymbol{z}) = p(\boldsymbol{x} \mid \boldsymbol{z})p(\boldsymbol{z})$$

p(z) = something simple $p(x \mid z) = g(z)$

Variational Auto-Encoder (VAE)

- Where does *z* come from? The classic DAG problem.
- The VAE approach: introduce an inference machine $q_{\phi}(z \mid x)$ that learns to approximate the posterior $p_{\theta}(z \mid x)$.
- Define a variational lower bound on the data likelihood: $p_{\theta}(x) \ge \mathcal{L}(\theta, \phi, x)$

$$\mathcal{L}(\theta, \phi, x) = \mathbb{E}_{q_{\phi}(z|x)} [\log z = \mathbb{E}_{q_{\phi}(z|x)} [\log z = -D_{\mathrm{KL}} (q_{\phi}(z|x))]$$

• What is $q_{\phi}(z \mid x)$?

 $g p_{\theta}(x, z) - \log q_{\phi}(z \mid x)]$ $g p_{\theta}(x \mid z) + \log p_{\theta}(z) - \log q_{\phi}(z \mid x)]$ $z \mid x) \parallel p_{\theta}(z)) + \mathbb{E}_{q_{\phi}(z \mid x)} \left[\log p_{\theta}(x \mid z)\right]$

zation term

reconstruction term



23

VAE Inference mode

lower bound:

$$\mathcal{L}(\theta, \phi, x) = -D_{\mathrm{KL}} \left(q_{\phi}(z \mid x) \| p_{\theta}(z) \right) + \mathbb{E}_{q_{\phi}(z \mid x)} \left[\log p_{\theta}(x \mid z) \right]$$

We parameterize $q_{\phi}(z \mid x)$ with another neural network:





The VAE approach: introduce an inference model $q_{\phi}(z \mid x)$ that learns to approximates the intractable posterior $p_{\theta}(z \mid x)$ by optimizing the variational





Reparametrization trick

- Adding a few details + one really important trick
- Let's consider z to be real and q_{ϕ}
- Parametrize z as $z = \mu_z(x) + c$
- (optional) Parametrize x as x =





$$\begin{aligned} (z \mid x) &= \mathcal{N}(z; \mu_z(x), \sigma_z(x)) \\ \sigma_z(x) \epsilon_z & \text{where } \epsilon_z = \mathcal{N}(0, 1) \\ \mu_x(z) &+ \sigma_x(z) \epsilon_x \text{ where } \epsilon_x = \mathcal{N}(0, 1) \end{aligned}$$

Training with backpropagation!

bound using gradient backpropagation.





• Due to a reparametrization trick, we can simultaneously train both the generative model $p_{\theta}(x \mid z)$ and the inference model $q_{\phi}(z \mid x)$ by optimizing the variational

Objective function: $\mathcal{L}(\theta, \phi, x) = -D_{\mathrm{KL}} \left(q_{\phi}(z \mid x) \| p_{\theta}(z) \right) + \mathbb{E}_{q_{\phi}(z \mid x)} \left[\log p_{\theta}(x \mid z) \right]$

vanilla VAE samples



Labelled Faces in the Wild (LFW)



ImageNet (small)



deep encoder/decoder: Some component collapse



Figures from Laurent Dinh & Vincent Dumoulin



deeper encoder/decoder: more component collapse

80



Figures from Laurent Dinh & Vincent Dumoulin

570 32232 $4 \frac{7}{2} \frac{1}{2} \frac{2}{3} \frac{6}{3} \frac{6}{3} \frac{1}{2} \frac{3}{2} \frac{$





PixelVAE



Ishaan Gulrajani, Kundan Kumar, Faruk Ahmed Adrien Ali Taiga, Francesco Visin, David Vazquez, Aaron Courville. ICLR 2017



Uses a PixelCNN in the VAE decoder to help avoid the blurring caused by the standard VAE assumption of independent pixels.

PixelVAE Samples (Gulrajani et al. 2017)



LSUN bedroom scenes (64x64)



ImageNet (64x64)





varying only the top-level latent variables

varying only the bottomlevel latent variables



varying only the pixel-level noise



















Inverse Autoregressive Flow (Kingma et al., NIPS 2016)



- the prior.
- much better fit between the posteriors and the prior.



(b) Posteriors in standard VAE (c) Posteriors in VAE with IAF

Standard VAE posteriors are factorized - limiting how well they can (marginally) fit

IAF greatly improves the flexibility of the posterior distributions, and allows for a

Normalizing Flows (Rezende and Mohamed, 2015)

- **Normalizing flows:** the transformation of a probability density through a sequence of invertible mappings.
- Transformation of random var For invertible functions:

$$q(\mathbf{z}') = q(\mathbf{z}) \left| \det \frac{\partial f^{-1}}{\partial \mathbf{z}'} \right| = q(\mathbf{z}) \left| \det \frac{\partial f}{\partial \mathbf{z}} \right|^{-1}$$
 by the Inverse Function Theorem

Chaining together a sequence

 $\log q_K(\boldsymbol{z}_K) = \log q_0$



iables:
$$oldsymbol{z}'=oldsymbol{f}(oldsymbol{z})$$
 , $oldsymbol{f}^{-1}(oldsymbol{z}')=oldsymbol{z}$

$$e: \boldsymbol{z}_{K} = \boldsymbol{f}_{K} \circ \boldsymbol{f}_{K-1} \circ \cdots \circ \boldsymbol{f}_{2} \circ \boldsymbol{f}_{1}(\boldsymbol{z}_{0}) \\ o(\boldsymbol{z}_{0}) - \sum_{k=1}^{K} \log \left| \det \frac{\partial \boldsymbol{f}_{k}}{\partial \boldsymbol{z}_{k}} \right|$$



Normalizing Flows (Rezende and Mohamed, 2015)

Law of the unconscious statistician: expectations w.r.t. the transformed density $q_K(z_K)$ can be written as expectations w.r.t. the original $q_0(z_0)$. For $z_K = f_K \circ f_{K-1} \circ \cdots \circ f_2 \circ f_1(z_0)$,

 $\mathbb{E}_{q_K}\left[g(\boldsymbol{z}_K)\right] = \mathbb{E}_{q_0}\left[g(\boldsymbol{f}_K)\right]$

The variational lower bound:

$$\begin{split} \mathcal{L}(\theta, \phi, \boldsymbol{x}) &= \mathbb{E}_{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})} \left[\log p_{\theta}(\boldsymbol{x}, \boldsymbol{z}) - \log q_{\phi}(\boldsymbol{z} \mid \boldsymbol{x}) \right] \\ &= \mathbb{E}_{q_{K}(\boldsymbol{z}_{K})} \left[\log p(\boldsymbol{x}, \boldsymbol{z}_{K}) - \log q_{K}(\boldsymbol{z}_{K}) \right] \\ &= \mathbb{E}_{q_{0}(\boldsymbol{z}_{0})} \left[\log p(\boldsymbol{x}, \boldsymbol{z}_{K}) - \log q_{0}(\boldsymbol{z}_{0}) + \sum_{k=1}^{K} \log \left| \det \frac{\partial \boldsymbol{f}_{k}}{\partial \boldsymbol{z}_{k-1}} \right| \right] \end{split}$$



$$_{K}\circ \boldsymbol{f}_{K-1}\circ\cdots\circ \boldsymbol{f}_{2}\circ \boldsymbol{f}_{1}(\boldsymbol{z}_{0}))\Big]$$



Approximate Posterior with In





Approximate Posterior with Inverse Autoregressive Flow (IAF)




VAE-IAF: MNIST log likelihood

Table 1: Generative modeling results on the dynamically sampled binarized MNIST version used in previous publications (Burda et al., 2015). Shown are averages; the number between brackets are standard deviations across 5 optimization runs. The right column shows an importance sampled estimate of the marginal likelihood for each model with 128 samples. Best previous results are reproduced in the first segment: [1]: (Salimans et al., 2014) [2]: (Burda et al., 2015) [3]: (Kaae Sønderby et al., 2016) [4]: (Tran et al., 2015)

Model

Convolutional VAE + HVI [1] DLGM 2hl + IWAE [2] LVAE [3] DRAW + VGP [4]





VLB	$\log p(\mathbf{x}) \approx$	
-83.49	-81.94	
	-82.90	
	-81.74	
-79.88		

Pottom-Up ResNet Block	Top-Down ResNet Block	

Another way to train a latent variable model





inference



Generative Adversarial Networks





Generative Adversarial Networks







GAN Objective

generator G with the minimax objective:

$$\min_{G} \max_{D} \mathbb{E}_{\boldsymbol{x} \sim \mathbb{P}_{r}} [\log(D(\boldsymbol{x}))] + \mathbb{E}_{\tilde{\boldsymbol{x}} \sim \mathbb{P}_{g}} [\log(1 - D(\tilde{\boldsymbol{x}}))].$$

where:

- \mathbb{P}_r is the data distribution

 $\tilde{\boldsymbol{x}} = G(\boldsymbol{z})$

- the generator input z is sampled from some simple noise distribution, (e.g. uniform or Gaussian).



• Formally, express the game between discriminator D and

- \mathbb{P}_q is the model distribution implicitly defined by:

$$z), \quad \boldsymbol{z} \sim p(\boldsymbol{z})$$

41

GAN Theory

• Optimal (nonparametric) discriminator:

 $D^*(\boldsymbol{x}) =$

Jensen-Shannon divergence between \mathbb{P}_r and \mathbb{P}_q .

$$JS(\mathbb{P}_r || \mathbb{P}_g) = KL\left(\mathbb{P}_r \left\|\frac{\mathbb{P}_r + \mathbb{P}_g}{2}\right) + KL\left(\mathbb{P}_g \left\|\frac{\mathbb{P}_r + \mathbb{P}_g}{2}\right)\right)$$

where $KL(\mathbb{P}_r || \mathbb{P}_g) = \int \log\left(\frac{p_r(x)}{p_g(x)}\right) p_r(x) d\mu(x)$



$$\frac{p_r(\boldsymbol{x})}{p_r(\boldsymbol{x}) + p_g(\boldsymbol{x})}$$

Under an ideal discriminator, the generator minimizes the





GAN Theory ... in practice

- discriminator saturates.
- training objective:
 - $\max_{D} \mathbb{E}_{\boldsymbol{x} \sim \mathbb{P}_{r}} [\log(D(\boldsymbol{x}))]$
 - $\max_{G} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim \mathbb{P}_{o}} [\log(D(\tilde{\boldsymbol{x}}))].$
- the presence of a good discriminator.



• The minimax objective leads to vanishing gradients as the

• In practice, Goodfellow et al (2014) advocate the heuristic

] +
$$\mathbb{E}_{\tilde{\boldsymbol{x}} \sim \mathbb{P}_g} [\log(1 - D(\tilde{\boldsymbol{x}}))].$$

However, this modified loss function can still misbehave in



GAN samples



MNIST



















































CIFAR-10













LEGST-SCUCIES GAN Xudong Mao, Qing Li[†], Haoran Xie, Raymond Y.K. Lau and Zhen Wang, ArXiv, Feb. 2017







128x128 LSUN bedroom scenes



GAN ZOO

GAN—Generative Adversarial Networks 3D-GAN—Learning a Probabilistic Latent Space of Object Shapes via 3D Generative-Adversarial Modeling acGAN—Face Aging With Conditional Generative Adversarial Networks AC-GAN—Conditional Image Synthesis With Auxiliary Classifier GANs AdaGAN—AdaGAN: Boosting Generative Models AEGAN—Learning Inverse Mapping by Autoencoder based Generative Adversarial Nets AffGAN—Amortised MAP Inference for Image Super-resolution AL-CGAN—Learning to Generate Images of Outdoor Scenes from Attributes and Semantic Layouts ALI—Adversarially Learned Inference AMGAN—Generative Adversarial Nets with Labeled Data by Activation Maximization AnoGAN—Unsupervised Anomaly Detection with Generative Adversarial Networks to Guide Marker Discovery ArtGAN—ArtGAN: Artwork Synthesis with Conditional Categorial GANs b-GAN—b-GAN: Unified Framework of Generative Adversarial Networks Bayesian GAN—Deep and Hierarchical Implicit Models BEGAN—BEGAN: Boundary Equilibrium Generative Adversarial Networks BiGAN—Adversarial Feature Learning BS-GAN—Boundary-Seeking Generative Adversarial Networks CGAN—Conditional Generative Adversarial Nets CCGAN—Semi-Supervised Learning with Context-Conditional Generative Adversarial Networks CatGAN—Unsupervised and Semi-supervised Learning with Categorical Generative Adversarial Networks CoGAN—Coupled Generative Adversarial Networks Context-RNN-GAN—Contextual RNN-GANs for Abstract Reasoning Diagram Generation C-RNN-GAN—C-RNN-GAN: Continuous recurrent neural networks with adversarial training CS-GAN—Improving Neural Machine Translation with Conditional Sequence Generative Adversarial Nets CVAE-GAN—CVAE-GAN: Fine-Grained Image Generation through Asymmetric Training CycleGAN—Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks DTN—Unsupervised Cross-Domain Image Generation DCGAN—Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks DiscoGAN—Learning to Discover Cross-Domain Relations with Generative Adversarial Networks DR-GAN—Disentangled Representation Learning GAN for Pose-Invariant Face Recognition DualGAN—DualGAN: Unsupervised Dual Learning for Image-to-Image Translation EBGAN—Energy-based Generative Adversarial Network f-GAN—f-GAN: Training Generative Neural Samplers using Variational Divergence Minimization GAWWN—Learning What and Where to Draw GoGAN—Gang of GANs: Generative Adversarial Networks with Maximum Margin Ranking GP-GAN—GP-GAN: Towards Realistic High-Resolution Image Blending IAN—Neural Photo Editing with Introspective Adversarial Networks iGAN—Generative Visual Manipulation on the Natural Image Manifold IcGAN—Invertible Conditional GANs for image editing ID-CGAN- Image De-raining Using a Conditional Generative Adversarial Network Improved GAN—Improved Techniques for Training GANs InfoGAN—InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets LAGAN—Learning Particle Physics by Example: Location-Aware Generative Adversarial Networks for Physics Synthesis LAPGAN—Deep Generative Image Models using a Laplacian Pyramid of Adversarial Networks LR-GAN—LR-GAN: Layered Recursive Generative Adversarial Networks for Image Generation LSGAN—Least Squares Generative Adversarial Networks

Deep Hunt, blog by Avinash Hindupur



LS-GAN—Loss-Sensitive Generative Adversarial Networks on Lipschitz Densities MGAN—Precomputed Real-Time Texture Synthesis with Markovian Generative Adversarial Networks MAGAN—MAGAN: Margin Adaptation for Generative Adversarial Networks MAD-GAN—Multi-Agent Diverse Generative Adversarial Networks MalGAN—Generating Adversarial Malware Examples for Black-Box Attacks Based on GAN MaliGAN—Maximum-Likelihood Augmented Discrete Generative Adversarial Networks MARTA-GAN—Deep Unsupervised Representation Learning for Remote Sensing Images McGAN—McGan: Mean and Covariance Feature Matching GAN MDGAN—Mode Regularized Generative Adversarial Networks MedGAN—Generating Multi-label Discrete Electronic Health Records using Generative Adversarial Networks MIX+GAN—Generalization and Equilibrium in Generative Adversarial Nets (GANs) MPM-GAN—Message Passing Multi-Agent GANs MV-BiGAN—Multi-view Generative Adversarial Networks pix2pix—Image-to-Image Translation with Conditional Adversarial Networks PPGN—Plug & Play Generative Networks: Conditional Iterative Generation of Images in Latent Space PrGAN—3D Shape Induction from 2D Views of Multiple Objects RenderGAN—RenderGAN: Generating Realistic Labeled Data RTT-GAN—Recurrent Topic-Transition GAN for Visual Paragraph Generation SGAN—Stacked Generative Adversarial Networks SGAN—Texture Synthesis with Spatial Generative Adversarial Networks SAD-GAN—SAD-GAN: Synthetic Autonomous Driving using Generative Adversarial Networks SalGAN—SalGAN: Visual Saliency Prediction with Generative Adversarial Networks SEGAN—SEGAN: Speech Enhancement Generative Adversarial Network SeGAN—SeGAN: Segmenting and Generating the Invisible SeqGAN—SeqGAN: Sequence Generative Adversarial Nets with Policy Gradient SimGAN—Learning from Simulated and Unsupervised Images through Adversarial Training SketchGAN—Adversarial Training For Sketch Retrieval SL-GAN—Semi-Latent GAN: Learning to generate and modify facial images from attributes Softmax-GAN—Softmax GAN SRGAN—Photo-Realistic Single Image Super-Resolution Using a Generative Adversarial Network S²GAN—Generative Image Modeling using Style and Structure Adversarial Networks SSL-GAN—Semi-Supervised Learning with Context-Conditional Generative Adversarial Networks StackGAN—StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks TGAN—Temporal Generative Adversarial Nets TAC-GAN—TAC-GAN—Text Conditioned Auxiliary Classifier Generative Adversarial Network TP-GAN—Beyond Face Rotation: Global and Local Perception GAN for Photorealistic and Identity Preserving Frontal View Synthesis Triple-GAN—Triple Generative Adversarial Nets Unrolled GAN—Unrolled Generative Adversarial Networks VGAN—Generating Videos with Scene Dynamics VGAN—Generative Adversarial Networks as Variational Training of Energy Based Models VAE-GAN—Autoencoding beyond pixels using a learned similarity metric VariGAN—Multi-View Image Generation from a Single-View ViGAN—Image Generation and Editing with Variational Info Generative AdversarialNetworks WGAN—Wasserstein GAN WGAN-GP—Improved Training of Wasserstein GANs WGAN-GP—Improved Training of Wasserstein GANs WaterGAN—WaterGAN: Unsupervised Generative Network to Enable Real-time Color Correction of Monocular Underwater Images



An explo-GAN of papers

Cumulative number of GAN papers by year





from Deep Hunt, blog by Avinash Hindupu^{#7}



DCGAN samples (Radford, Metz and Chintala; 2016)

Z-space interpolations

LSUN bedroom scenes





Cartoon of the Image manifold



What makes GANs special?

 \mathcal{X}_1

more traditional max-likelihood approach

GAN

Training a GAN: Distances between Manifolds

- Data manifold
- GAN manifold (Generative model)

51

Training a GAN: Distances between Manifolds

Data manifold

Jensen-Shannon Divergence

$$\mathrm{JS}(\mathbb{P}_r \| \mathbb{P}_g) = \mathrm{KL}\left(\mathbb{P}_r \left\| \frac{\mathbb{P}_r + \mathbb{P}_g}{2} \right) + \mathrm{KL}\left(\mathbb{P}_g \left\| \frac{\mathbb{P}_r + \mathbb{P}_g}{2} \right\| \right)\right)$$

• What is the JS divergence in this simple case?

$$\mathrm{JS}(\mathbb{P}_r \| \mathbb{P}_g) = \begin{cases} \log 2 & \text{if } \theta \neq 0, \\ 0 & \text{if } \theta = 0, \end{cases}$$

Jensen-Shannon Divergence

Earth-Movers Distance

- distance.

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

- distribution \mathbb{P}_r into the distribution \mathbb{P}_q .

• JS divergence is not a useful learning signal to train GANs.

• An alternative: Earth-Mover (also called Wasserstein-1)

Minimum cost of transporting mass to transform the

The EM distance is continuous everywhere and differentiable almost everywhere (under mild assumptions).

Wasserstein Distance

- $W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} \left[\|x y\| \right]$
 - What is the EM (or Wasserstein) distance in this simple case?

Wasserstein Distance

 $W(\mathbb{P}_r \| \mathbb{P}_g) = |\theta|$

Wasserstein Distance

 $W(\mathbb{P}_r \| \mathbb{P}_g) = |\theta|$

Wasserstein GAN Arjovsky, Chintala, Bottou (2017)

- $W(\mathbb{P}_r, \mathbb{P}_q)$ might have nice properties compared to $JS(\mathbb{P}_r, \mathbb{P}_q)$
- However, the infimum is intractable in:

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(f)} V(g) = \sum_{\gamma \in \Pi(f)} V(g) = \sum_{\Pi$$

Can exploit Kantorovich-Rubinstein duality:

$$W(\mathbb{P}_r, \mathbb{P}_g) = \sup_{\|f\|_L \le 1} \mathbb{E}_{x \sim \mathbb{P}_r}[f(x)] - \mathbb{E}_{x \sim \mathbb{P}_g}[f(x)]$$

where the supremum is over all the 1-Lipschitz functions $f: \mathcal{X} \to \mathbb{R}$

 $\inf_{(\mathbb{P}_r,\mathbb{P}_a)} \mathbb{E}_{(x,y)\sim\gamma} \left[\|x-y\| \right]$

Wasserstein GAN Arjovsky, Chintala, Bottou (2017)

- The WGAN Objective function:
 - $\min_{G} \max_{D \in \mathcal{D}} \mathbb{E} \left[L \\ \mathbf{x} \sim \mathbb{P}_r \right]^r$
 - where \mathcal{D} is the set of 1-Lipschitz functions.
- on the critic *D*?
 - lie within a compact space [-c, c].
 - —

$$D(oldsymbol{x})ig] - \mathop{\mathbb{E}}\limits_{ ilde{oldsymbol{x}} \sim \mathbb{P}_g} ig[D(ilde{oldsymbol{x}}))ig]$$

• Open question: how to effectively enforce the Lipschitz constraint

- Arjovsky et al. (2017) propose to clip the weights of the critic to

Results in a subset of the k-Lipschitz functions (k is a function of c).

Issues with Weight Clipping

1. Underuse capacity 2. Exploding and vanishing gradients

Gradient Penalty Approach Gulrajani, Ahmed, Arjovsky, Dumoulin, Courville (2017)

- line between $m{x}$ and $m{ ilde{x}}$) then: $\nabla D^*(\boldsymbol{x}_t)$
- \bullet

$$L = \underbrace{\mathbb{E}}_{\tilde{\boldsymbol{x}} \sim \mathbb{P}_g} \left[D(\tilde{\boldsymbol{x}}) \right] - \underbrace{\mathbb{E}}_{\boldsymbol{x} \sim \mathbb{P}_r} \left[D(\boldsymbol{x}) \right]_{r \sim \mathbb{P}_r}$$

• A property of the optimal WGAN critic: If $ilde{x} \sim \mathbb{P}_q$ then there is a point $m{x} \sim \mathbb{P}_r,$ such that for all points $m{x}_t = tm{x} + (1 - t) ilde{m{x}}$ (on a straight

$$)=rac{oldsymbol{x}-oldsymbol{x}_t}{\|oldsymbol{x}-oldsymbol{x}_t\|}$$

This implies the optimal WGAN critic has gradient norm 1 at x_t

Gradient Penalty version of WGAN (i.e. the WGAN-GP) objective: $(\boldsymbol{x})] + \lambda \mathop{\mathbb{E}}_{\hat{\boldsymbol{x}} \sim \mathbb{P}_{\hat{\boldsymbol{x}}}} \left[(\|\nabla_{\hat{\boldsymbol{x}}} D(\hat{\boldsymbol{x}})\|_2 - 1)^2 \right]$

Our gradient penalty

Gulrajani, Ahmed, Arjovsky, Dumoulin, Courville (2017)

Gradient penalty:

$$\mathbb{E}_{\hat{\boldsymbol{x}} \sim \mathbb{P}_{\hat{\boldsymbol{x}}}} \left[(\|\nabla_{\hat{\boldsymbol{x}}} D(\hat{\boldsymbol{x}})\|_2 - 1)^2 \right]$$

Sample along straight lines:

$$\epsilon \sim U[0, 1], \boldsymbol{x} \sim \mathbb{P}_r, \tilde{\boldsymbol{x}} \sim \mathbb{P}_g$$

 $\hat{\boldsymbol{x}} = \epsilon \boldsymbol{x} + (1 - \epsilon) \tilde{\boldsymbol{x}}$

Comparison on difficult to train architectures

DCGAN

- Comparison based on recommended default parameter setting for each algorithm.
- WGAN-GP is more robust to variations in training setups.

Baseline (G: DCGAN, D: DCGAN)

WGAN (clipping)

LSGAN

WGAN-GP (ours)

G: No BN and a constant number of filters, D: DCGAN

G: 4-layer 512-dim ReLU MLP, D: DCGAN

Comparison on difficult to train architectures

DCGAN

- Comparison based on recommended default parameter setting for each algorithm.
- WGAN-GP is more robust to variations in training setups.

Gated multiplicative nonlinearities everywhere in G and D

101-layer ResNet G and D

WGAN (clipping)

WGAN-GP (ours)

tanh nonlinearities everywhere in G and D

WGAN with Gradient Penalty

Convergence on CIFAR-10

But what about inference...

- How can we use generative models?
 - GANs can generate content, but somethings you want to make inference about observed data.
- Can we incorporate an inference mechanism into GANs?
- Can we learn an inference mechanisms using an adversarial training paradigm?

Iwo papers, one mode

- LEARNED INFERENCE, arXiv:1606.00704
- FEATURE LEARNING, arXiv:1605.09782

• **ALI**: Vincent Dumoulin, Ishmael Belghazi, Olivier Mastropietro Ben Poole, Alex Lamb, Martin Arjovsky (2016) ADVERSARIALLY

• **BiGAN**: Donahue, Krähenbühl and Darrell (2016), ADVERSARIAL

• But also showing results on Hierarchical ALI by Ishmael Belghazi, Sai Rajeshwar, Olivier Mastropietro and Negar Rostamzadeh

Adversarially learned inference: Main idea

- Cast the learning of both an inference model (encoder) and a generative model (*decoder*) in a GAN-like adversarial framework.
- Discriminator is trained to discriminate between *joint* samples (**x**, **z**) \bullet from:
 - Data distribution - Encoder distribution $q(\mathbf{x}, \mathbf{z}) = q(\mathbf{x})q(\mathbf{z} | \mathbf{x})$, or - Decoder distribution $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{z})p(\mathbf{x} | \mathbf{z})$. **Prior distribution**
- Generator learns conditionals $q(\mathbf{z} \mid \mathbf{x})$ and $p(\mathbf{x} \mid \mathbf{z})$ to fool the discriminator.

ALL: model diagram

70

Toy Example

• Learning the Identity function:

Encoder: $X \sim N(0,1)$ Decoder: $Z \sim N(0,1)$

Zihang Dai

----- Update 0 ------

Theoretical properties

In analogy with GAN, under an ideal the Jensen-Shannon divergence between $p(\mathbf{x}, \mathbf{z})$ and $q(\mathbf{x}, \mathbf{z})$.

discriminator, the generator minimizes

72




Samples





CelebA face dataset



Samples





Hierarchical ALI: model diagram







Hierarchical ALI





CelebA-128X128



















Model samples





































































Reconstructions given z_1 , z_2

Reconstructions given z₂

ImageNet-128X128

Model samples

Hierarchical ALI

Unconditional

























Reconstructions given z₁, z₂

Recon

Data



Reconstructions given z₂



Generative models II: Outline

- Autoregressive models
 - PixelCNN
- Latent variable models
 - Variational Autoencoders
 - VAEs
 - Inverse Autoregressive Flow: An improved strategy for inference
 - Generative Adversarial Networks
 - GAN
 - Wasserstein GAN
 - ALI



