# Generative Adversarial Networks (GANs)

Ian Goodfellow, OpenAI Research Scientist Presentation at Berkeley Artificial Intelligence Lab, 2016-08-31



### Generative Modeling

• Density estimation



• Sample generation



Training examples

(Goodfellow 2016)

#### Maximum Likelihood





### Fully Visible Belief Nets

• Explicit formula based on chain (Frey et al, 1996) rule:  $p_{\text{model}}(\boldsymbol{x}) = p_{\text{model}}(x_1) \prod_{i=1}^{n} p_{\text{model}}(x_i \mid x_1, \dots, x_{i-1})$ 

i=2

- Disadvantages:
  - O(n) sample generation cost
  - Currently, do not learn a useful latent representation



PixelCNN elephants (van den Ord et al 2016)

### Change of Variables

$$y = g(x) \Rightarrow p_x(x) = p_y(g(x)) \left| \det \left( \frac{\partial g(x)}{\partial x} \right) \right|$$

#### e.g. Nonlinear ICA (Hyvärinen 1999)

Disadvantages:



64x64 ImageNet Samples Real NVP (Dinh et al 2016)

- Transformation must be invertible
- Latent dimension must match visible dimension

#### Variational Autoencoder (Kingma and Welling 2013, Rezende et al 2014)

 $\log p(\boldsymbol{x}) \ge \log p(\boldsymbol{x}) - D_{\mathrm{KL}} \left( q(\boldsymbol{z}) \| p(\boldsymbol{z} \mid \boldsymbol{x}) \right)$  $= \mathbb{E}_{\boldsymbol{z} \sim q} \log p(\boldsymbol{x}, \boldsymbol{z}) + H(q)$ 





Disadvantages: -Not asymptotically consistent unless q is perfect -Samples tend to have lower quality

CIFAR-10 samples (Kingma et al 2016)

# Boltzmann Machines $p(\boldsymbol{x}) = \frac{1}{Z} \exp(-E(\boldsymbol{x}, \boldsymbol{z}))$ $Z = \sum_{\boldsymbol{x}} \sum_{\boldsymbol{z}} \exp(-E(\boldsymbol{x}, \boldsymbol{z}))$

- Partition function is intractable
- May be estimated with Markov chain methods
- Generating samples requires Markov chains too

#### GANs

- Use a latent code
- Asymptotically consistent (unlike variational methods)
- No Markov chains needed
- Often regarded as producing the best samples
  - No good way to quantify this

## Generator Network $\boldsymbol{x} = G(\boldsymbol{z}; \boldsymbol{\theta}^{(G)})$



- -Must be differentiable
- In theory, could use REINFORCE for discrete variables
- No invertibility requirement
- Trainable for any size of z
- Some guarantees require z to have higher dimension than x
- Can make x conditionally Gaussian given z but need not do so

### Training Procedure

- Use SGD-like algorithm of choice (Adam) on two minibatches simultaneously:
  - A minibatch of training examples
  - A minibatch of generated samples
- Optional: run k steps of one player for every step of the other player.

#### Minimax Game

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}} \log D(\boldsymbol{x}) - \frac{1}{2} \mathbb{E}_{\boldsymbol{z}} \log \left(1 - D\left(G(\boldsymbol{z})\right)\right)$$
$$J^{(G)} = -J^{(D)}$$

- -Equilibrium is a saddle point of the discriminator loss
- -Resembles Jensen-Shannon divergence
- -Generator minimizes the log-probability of the discriminator being correct

Non-Saturating Game  

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}} \log D(\boldsymbol{x}) - \frac{1}{2} \mathbb{E}_{\boldsymbol{z}} \log (1 - D(G(\boldsymbol{z})))$$

$$J^{(G)} = -\frac{1}{2} \mathbb{E}_{\boldsymbol{z}} \log D(G(\boldsymbol{z}))$$

-Equilibrium no longer describable with a single loss -Generator maximizes the log-probability of the discriminator being mistaken

-Heuristically motivated; generator can still learn even when discriminator successfully rejects all generator samples

#### Maximum Likelihood Game

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}} \log D(\boldsymbol{x}) - \frac{1}{2} \mathbb{E}_{\boldsymbol{z}} \log \left(1 - D\left(G(\boldsymbol{z})\right)\right)$$
$$J^{(G)} = -\frac{1}{2} \mathbb{E}_{\boldsymbol{z}} \exp\left(\sigma^{-1}\left(D\left(G(\boldsymbol{z})\right)\right)\right)$$

When discriminator is optimal, the generator gradient matches that of maximum likelihood

("On Distinguishability Criteria for Estimating Generative Models", Goodfellow 2014, pg 5)

#### Maximum Likelihood Samples



### Discriminator Strategy

Optimal  $D(\boldsymbol{x})$  for any  $p_{\text{data}}(\boldsymbol{x})$  and  $p_{\text{model}}(\boldsymbol{x})$  is always

A cooperative rather than Discriminator adversarial view of GANs: the discriminator tries to estimate the ratio of the data and model distributions, and informs the generator of its estimate in order to guide its improvements.



#### Comparison of Generator Losses



(Goodfellow 2016)

#### DCGAN Architecture



(Radford et al 2015)

#### DCGANs for LSUN Bedrooms



#### (Radford et al 2015)

### Vector Space Arithmetic





Man



Man with glasses Woman



Woman with Glasses

### Mode Collapse

- Fully optimizing the discriminator with the generator held constant is safe
- Fully optimizing the generator with the discriminator held constant results in mapping all points to the argmax of the discriminator
- Can partially fix this by adding nearest-neighbor features constructed from the current minibatch to the discriminator ("minibatch GAN") (Salimans et al 2016)

#### Minibatch GAN on CIFAR



Training Data

Samples

(Salimans et al 2016)

#### Minibatch GAN on ImageNet



(Salimans et al 2016)

### Cherry-Picked Results















#### GANs Work Best When Output Entropy is Low

this small bird has a pink breast and crown, and black almost all black with a red primaries and secondaries.



the flower has petals that are bright pinkish purple with white stigma

this magnificent fellow is crest, and white cheek patch.



this white and yellow flower have thin white petals and a round yellow stamen





(Reed et al 2016)

### Optimization and Games

Optimization: find a minimum:

$$\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

Game:

Player 1 controls  $\boldsymbol{\theta}^{(1)}$ Player 2 controls  $\boldsymbol{\theta}^{(2)}$ Player 1 wants to minimize  $J^{(1)}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)})$ Player 2 wants to minimize  $J^{(2)}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)})$ Depending on J functions, they may compete or cooperate.

#### Games $\supseteq$ optimization

Example:  $\theta^{(1)} = \theta$   $\theta^{(2)} = \{\}$   $J^{(1)}(\theta^{(1)}, \theta^{(2)}) = J(\theta^{(1)})$  $J^{(2)}(\theta^{(1)}, \theta^{(2)}) = 0$ 

### Nash Equilibrium

• No player can reduce their cost by changing their own strategy:

$$\forall \boldsymbol{\theta}^{(1)}, J^{(1)}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)*}) \ge J^{(1)}(\boldsymbol{\theta}^{(1)*}, \boldsymbol{\theta}^{(2)*})$$
  
$$\forall \boldsymbol{\theta}^{(2)}, J^{(2)}(\boldsymbol{\theta}^{(1)*}, \boldsymbol{\theta}^{(2)}) \ge J^{(2)}(\boldsymbol{\theta}^{(1)*}, \boldsymbol{\theta}^{(2)*})$$

- In other words, each player's cost is minimal with respect to that player's strategy
- Finding Nash equilibria  $\subseteq$  optimization (but not clearly useful)

#### Well-Studied Cases

- Finite minimax (zero-sum games)
- Finite mixed strategy games
- Continuous, convex games
- Differential games (lion chases gladiator)



#### Local Differential Nash Equilibria

$$\nabla_{\boldsymbol{\theta}^{(i)}} J^{(i)}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}) = \mathbf{0}$$

Necessary:

 $\nabla^2_{\boldsymbol{\theta}^{(i)}} J^{(i)}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)})$  is positive semi-definite Sufficient:

 $\nabla^2_{\boldsymbol{\theta}^{(i)}} J^{(i)}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)})$  is positive definite

(Ratliff et al 2013)

#### Sufficient Condition for Simultaneous Gradient Descent to Converge

$$\boldsymbol{\omega} = \left[ \begin{array}{c} \nabla_{\boldsymbol{\theta}^{(1)}} J^{(1)}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}) \\ \nabla_{\boldsymbol{\theta}^{(2)}} J^{(2)}(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}) \end{array} \right]$$

The eigenvalues of  $\nabla_{\theta} \omega$  must have positive real part:

$$\begin{bmatrix} \nabla^2_{\theta^{(1)}} J^{(1)} & \nabla_{\theta^{(1)}} \nabla_{\theta^{(2)}} J^{(2)} \\ \nabla_{\theta^{(2)}} \nabla_{\theta^{(1)}} J^{(1)} & \nabla^2_{\theta^{(2)}} J^{(2)} \end{bmatrix}$$
(I call this the "generalized Hessian")

(Ratliff et al 2013)

### Interpretation

- Each player's Hessian should have large, positive eigenvalues, expressing a strong preference to keep doing their current strategy
- The Jacobian of one player's gradient with respect to the other player's parameters should have smaller contributions to the eigenvalues, meaning each player has limited ability to change the other player's behavior at convergence
- Does not apply to GANs, so their convergence remains an open question

#### Equilibrium Finding Heuristics

- Keep parameters near their running average
  - Periodically assign running average value to parameters
  - Constrain parameters to lie near running average
  - Add loss for deviation from running average

#### Stabilized Training



### Other Games in AI

- Robust optimization / robust control
  - for security/safety, e.g. resisting adversarial examples
- Domain-adversarial learning for domain adaptation
- Adversarial privacy
- Guided cost learning
- Predictability minimization

#### Conclusion

- GANs are generative models that use supervised learning to approximate an intractable cost function
- GANs can simulate many cost functions, including the one used for maximum likelihood
- Finding Nash equilibria in high-dimensional, continuous, non-convex games is an important open research problem