

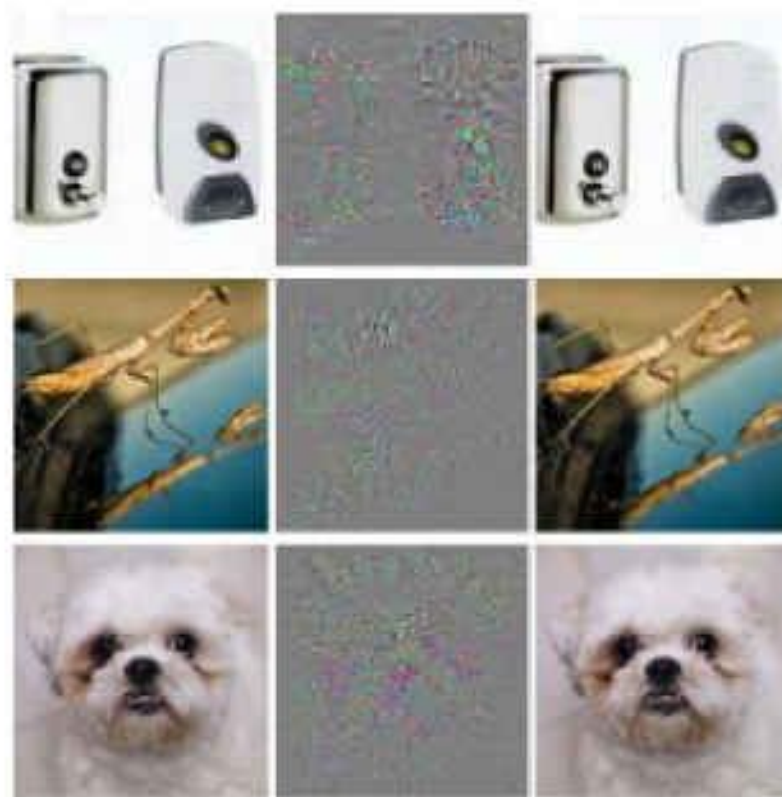
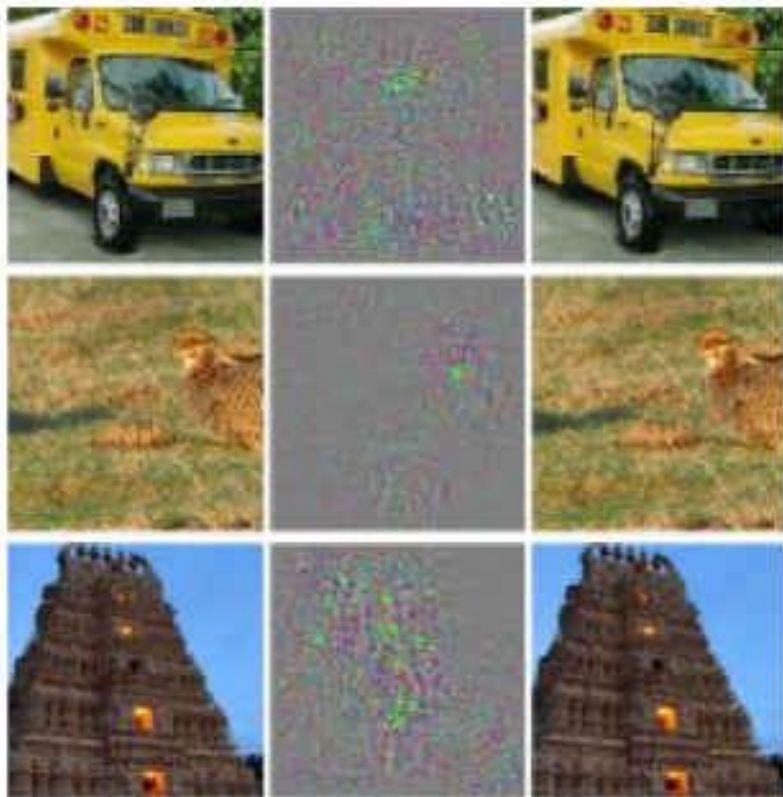
Generative Adversarial Networks

The success of deep learning thus far

- So far, we've been talking about deep learning in the context of classification
- Today we will see that
 - They aren't as good at classification as it might seem
 - They can't generate images at all

Deep Neural Networks are Easily Fooled

Left: correctly classified; Right: Incorrectly classified



Generative Adversarial Networks

- We are going to address both of these in one algorithm:
 - Learn to generate images
 - Stop being susceptible to adversarial examples

Generative Adversarial Networks

- We will define two networks
 - $G(z; \theta_g)$ is a generator
 - Takes input noise z and has parameters θ_g
 - $D(x; \theta_d)$ is a pre-trained discriminator, 0-1 probability
 - Takes input x and has parameters θ_d
- We train D to *maximize* probability of correctly classifying training examples vs. data generated by G
- We train G to *minimize* $\log(1 - D(G(z)))$

Generative Adversarial Networks

- We train D to *maximize* probability of correctly classifying training examples vs. data generated by G ; i.e. $D(G(z))$
- We train G to *minimize* $\log(1 - D(G(z)))$
- So $D(G(z))$ is trying to reach 1
- But $\log(1 - D(G(z)))$ is trying to reach $-\infty$
 - In other words, the two networks play a minimax game with value function $V(D, G)$. This is just an version of continuous cross-entropy loss with two mini-batches rather than 1.

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))] .$$

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**

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of m examples $\{x^{(1)}, \dots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- Update the discriminator by ascending its stochastic gradient:

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

end for

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Update the generator by descending its stochastic gradient:

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

end for

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

4.1 Global Optimality of $p_g = p_{\text{data}}$

We first consider the optimal discriminator D for any given generator G .

Proposition 1. *For G fixed, the optimal discriminator D is*

$$D_G^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})} \quad (2)$$

Proof. The training criterion for the discriminator D , given any generator G , is to maximize the quantity $V(G, D)$

$$\begin{aligned} V(G, D) &= \int_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \log(D(\mathbf{x})) d\mathbf{x} + \int_{\mathbf{z}} p_z(\mathbf{z}) \log(1 - D(g(\mathbf{z}))) d\mathbf{z} \\ &= \int_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \log(D(\mathbf{x})) + p_g(\mathbf{x}) \log(1 - D(\mathbf{x})) d\mathbf{x} \end{aligned} \quad (3)$$

For any $(a, b) \in \mathbb{R}^2 \setminus \{0, 0\}$, the function $y \rightarrow a \log(y) + b \log(1 - y)$ achieves its maximum in $[0, 1]$ at $\frac{a}{a+b}$. The discriminator does not need to be defined outside of $\text{Supp}(p_{\text{data}}) \cup \text{Supp}(p_g)$, concluding the proof. \square

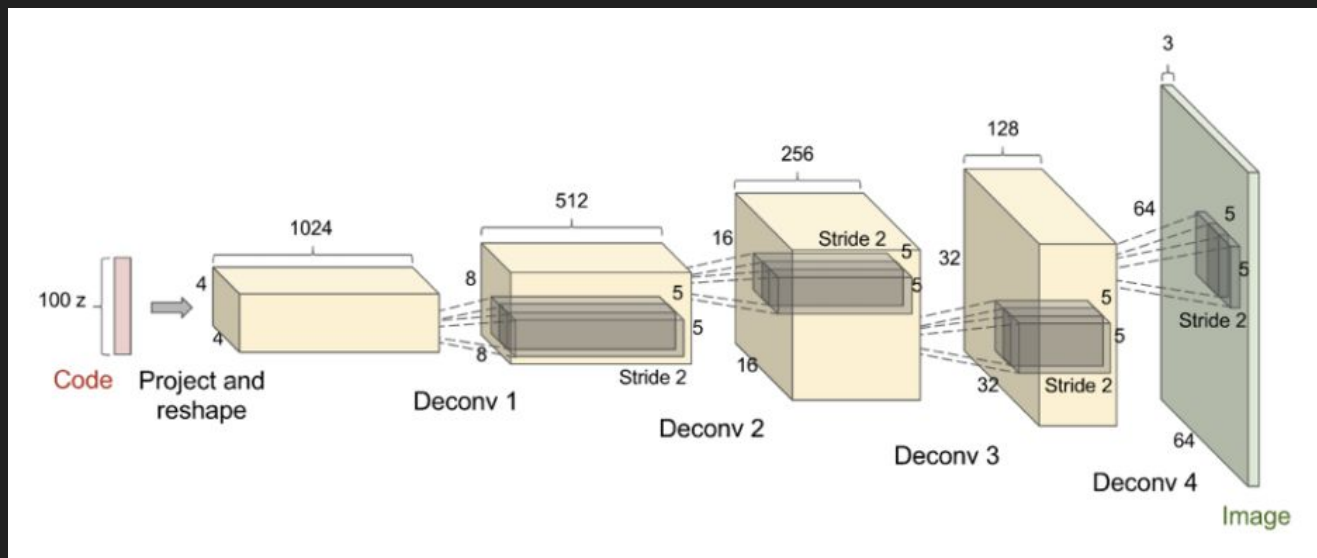
Note that the training objective for D can be interpreted as maximizing the log-likelihood for estimating the conditional probability $P(Y = y|\mathbf{x})$, where Y indicates whether \mathbf{x} comes from p_{data} (with $y = 1$) or from p_g (with $y = 0$). The minimax game in Eq. 1 can now be reformulated as:

Further Proofs

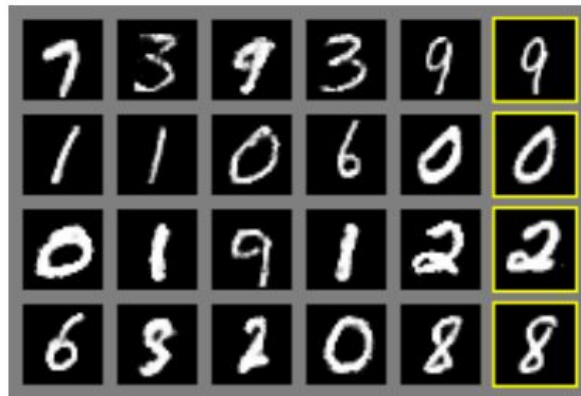
- They then go on to prove that given the algorithm, the generated data will converge to the true data samples
- Steps:
 - Show that the function $V(D, G)$ is convex
 - Show that its global optima is at the location where $p_g = p_{\text{data}}$
 - p_g is the probability distribution of samples $G(z)$
 - $z \sim p_z =$ probability distribution of random seeds

Model Architecture

- What are potential shortcomings of this architecture?



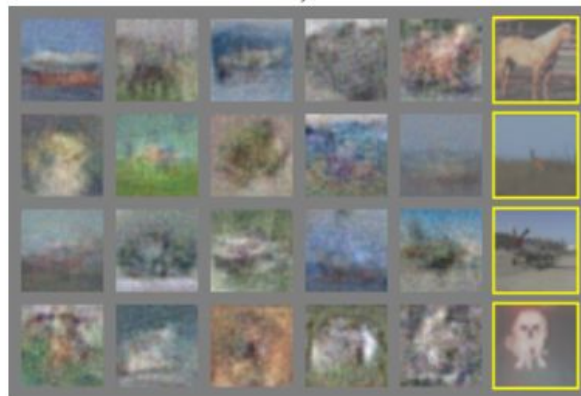
Results



a)



b)

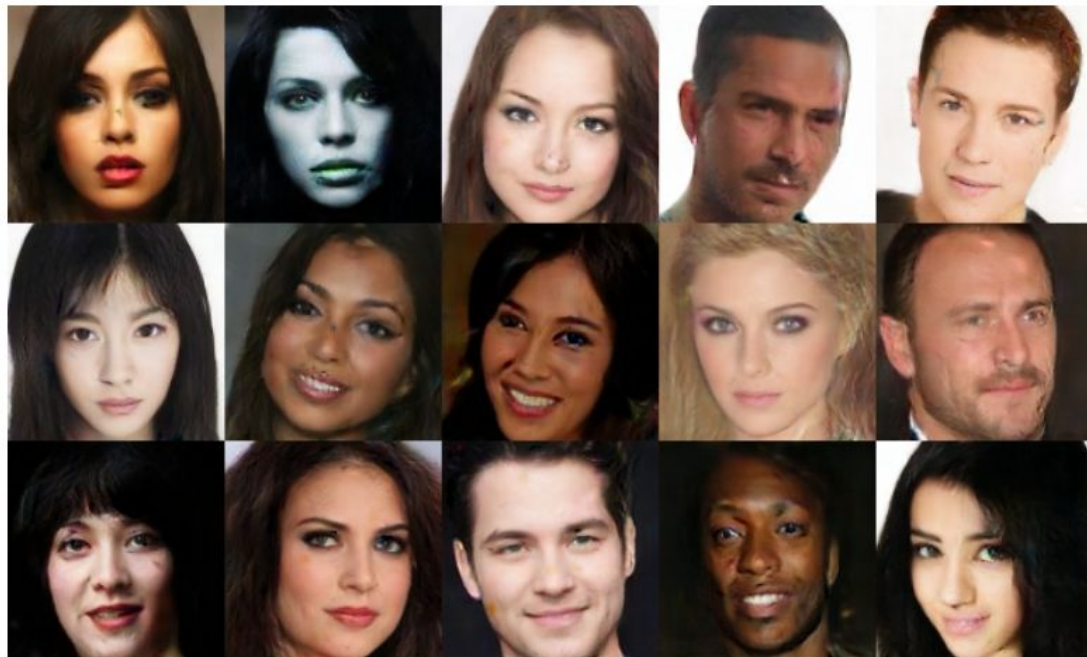


c)



d)

BE-GAN



(b) Our results (128x128)

Impact

- Discriminators are less susceptible to adversarial examples
- Generators are 'as good' as the true data samples
 - At least in the eyes of the discriminators
- Still not great at regular every day images

Further Reading
by Ian Goodfellow at OpenAI

What we've learned over the semester

- Machine Learning Basics
- Neural Network Math
 - Backpropagation and optimization
 - Batch normalization, dropout, L2 norm
- Neural Network Architectures
 - LSTM, ConvNet, Auto-Encoder
- 5 Difficult (!!!) Papers

If the material from this course makes sense...

- Then you're basically ready for a full-time job working in deep learning
- May need some practice writing the code and learning TensorFlow or Theano
- Next steps for deep learning proficiency
 - Read papers!!!!!! I have a GoogleDoc of some interesting ones I can share, just email me.
 - [Reddit.com/r/machinelearning](https://www.reddit.com/r/machinelearning)
 - Just look at this [thread](#) from today for example
 - Participate in Github projects, keep up with the latest code
 - Take linear algebra, probability, math-statistics, calculus