

# Announcements

① HW 2 now due Monday, 5/10

Student holiday on Fri!

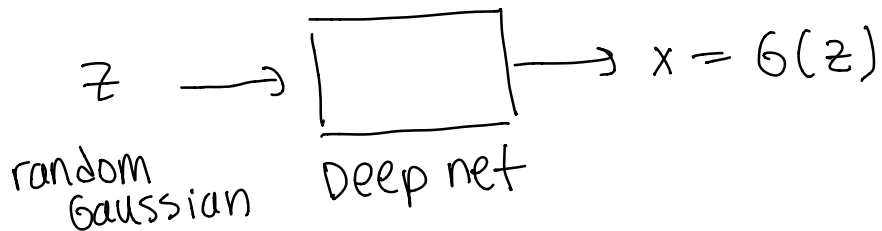
② Still doing project meetings...

# Deep Generative Models

Deep learning is useful for more than just supervised learning

Main Question: Can we learn to generate realistic new data?

The usual setup is:



**But what does realistic mean?**

Goal: No deep network can discriminate between outputs of  $G$  and real samples

Goodfellow et al. defined a min-max objective function

$$\min_u \max_v \mathbb{E}_{x \sim D_{\text{real}}} [\log D_v(x)] + \mathbb{E}_{x \sim G_u} [\log (1 - D_v(x))]$$

where:

- ①  $D_{\text{real}}$  is the distribution on real data
- ②  $D_{G_u}$  the distribution on outputs of generator  $G_u$  and  $u$  are its params
- ③  $D_v$  is the discriminator, that tries to output 1 on real data, 0 o.w

Then we try to find an equilibrium, i.e.

"A generator  $G_u$  so that no discrim can beat random guessing."

Some things to worry about

① How do we solve this optimization problem?

As usual, we use SGD or some other heuristic

② what if we only have a finite (polynomial) sized set of samples?

③ what if instead of converging, we cycling?

Let's come back to these later, and first study the properties at equilibrium

Lemma 1 If  $D_r$  is allowed to be any function from  $\mathbb{R}^d$  to  $[0,1]$  then the optimal choice is:

$$D_V(x) = \frac{P_{\text{real}}(x)}{P_{\text{real}}(x) + P_{\text{fake}}(x)}$$

where  $P_{\text{real}}$  and  $P_{\text{fake}}$  are the density functions of  $D_{\text{real}}$  and  $D_{\text{fake}}$  respectively

In particular, if the supports of  $D_{\text{real}}$  and  $D_{\text{fake}}$  are disjoint, then the optimal  $D_V$  achieves zero loss

The intuition behind the proof is simple.  
Let

$H_0$ :  $X$  came from  $D_{\text{real}}$

$H_1$ :  $X$  came from  $D_{\text{fake}}$

Then the optimal  $D_V$  computes the optimal hypothesis test

Our actual proof will use the KL-divergence

def: The KL-divergence between discrete distributions  $p$  and  $q$  is

$$D_{KL}(p||q) \triangleq \sum_x p(x) \log \frac{p(x)}{q(x)}$$

Or if  $p$  and  $q$  are continuous

$$D_{KL}(p||q) = \int p(x) \log \frac{p(x)}{q(x)} dx$$

We will use the following elementary fact

Fact 1: For distributions  $p$  and  $q$ ,

$D_{KL}(p||q) \geq 0$  and equality is achieved iff  $p=q$  almost surely

Rearranging things, we get

Corollary 1: For distributions  $p$  and  $q$ ,

$\sum_x p(x) \log q(x)$  is maximized for  $p=q$

Now we can prove the lemma

Proof: We can rewrite the objective function in (\*) as

$$\int (P_{\text{real}}(x) \log D_r(x) + P_{G_u}(x) \log(1 - D_r(x))) dx$$

and we can maximize this (over  $D_r$ ) on an  $x$ -by- $x$  basis, and from Corollary 1 the optimal choice is

$$D_r(x) = \frac{P_{\text{real}}(x)}{P_{\text{real}}(x) + P_{G_u}(x)}$$

as desired. 

Now let's plug in the optimal  $D_r$  and see what we get for the optimal choice of the generator  $D_{G_u}$

$$\underbrace{\int P_{\text{real}}(x) \log \frac{P_{\text{real}}(x)}{P_{\text{real}}(x) + P_{G_u}(x)} dx}_{\text{KL Divergence}} + \underbrace{\int P_{G_u}(x) \log \frac{P_{G_u}(x)}{P_{\text{real}}(x) + P_{G_u}(x)} dx}_{\text{KL Divergence}}$$

$$= D_{\text{KL}}\left(P_{\text{real}} \parallel \frac{P_{\text{real}} + P_{G_u}}{2}\right) + D_{\text{KL}}\left(P_{G_u} \parallel \frac{P_{\text{real}} + P_{G_u}}{2}\right)$$

def: For two distributions  $p$  and  $q$ , the Jensen-Shannon divergence is

$$D_{\text{JS}}(p \parallel q) \triangleq \frac{1}{2} D_{\text{KL}}\left(p \parallel \frac{p+q}{2}\right) + \frac{1}{2} D_{\text{KL}}\left(q \parallel \frac{p+q}{2}\right)$$

Theorem [Goodfellow et al] If  $D_V$  is chosen among all functions  $\mathbb{R}^d \rightarrow [0, 1]$ , then at equilibrium the optimal generator  $m$  minimizes  $D_{\text{JS}}(D_{\text{real}} \parallel D_{G_u})$

Training GANs following (\*) is tricky



Main Question: What happens if we change the game? Characterize equilibrium?

Arjovsky et al. first considered:

$$\min_u \max_v \mathbb{E}_{x \sim D_{\text{real}}} [D_v(x)] + \mathbb{E}_{x \sim D_{G_u}} [1 - D_v(x)] \quad (\square)$$

To understand its equilibrium, we will need yet another distance on distributions

def: The Wasserstein distance between  $p$  and  $q$  is

$$W_1(p, q) = \sup_{D \text{ is 1-lipschitz}} \left| \mathbb{E}_{x \sim p} [D(x)] - \mathbb{E}_{x \sim q} [D(x)] \right|$$

Let's see some examples to get some intuition

Suppose  $p$  and  $q$  are discrete distributions



the Wasserstein distance measures

$$\text{cost} \triangleq \text{mass} \times \text{distance}$$

of moving  $p$  into  $q$

Thus it is a smoother distance, e.g. it distinguishes between



What can we say about equilibrium?

Theorem 2 [Arjovsky et al.] If  $D_V$  is chosen optimally all 1-Lipschitz functions, then optimal generator minimizes

$$W_1(D_{\text{real}}, D_G)$$

Now let's dig into some of the theoretical issues

Issue #1: We can only estimate

$$\mathbb{E}_{x \sim D_{\text{real}}} [\log(D_V(x))], D_{\text{JS}}(D_G \| D_{\text{real}}), W_1(P, q)$$

from a polynomial # of samples

i.e. does  $W_1(P, q)$  being small imply that  $W_1(\hat{P}, \hat{q})$  is too

↑ ↑  
empirical distributions

Definitely not!

Observation 1 [Arora et al.] Consider  $P = \mathcal{N}(0, \frac{1}{2}I)$ . Then whp for poly. samples

$$d_{JS}(P \parallel \hat{P}) = \log 2 \quad (\text{not prob.})$$

$$W_1(P \parallel \hat{P}) \geq 1.1$$

Various works have proposed using the Sliced Wasserstein distance instead

Arora, Ge, Liang, Ma, Zhang define a distance based on neural networks:

def. Let  $\mathcal{F}$  be a class of neural nets with  $n$  parameters. Let  $\phi$  be concave

$$D_{\mathcal{F}, \phi}(P, Q) \triangleq \sup_{D \in \mathcal{F}} \mathbb{E}_{x \sim P} [\phi(D(x))] + \mathbb{E}_{x \sim Q} [\phi(1-D(x))] - 2\phi\left(\frac{1}{2}\right)$$

They study generalization and (approximate) equilibria

Theorem 3 [Arora et al]. Suppose that

- ①  $\phi$  is concave,  $L_\phi$ -Lipschitz and takes values in  $[-\Delta, \Delta]$
- ② The class of discriminators is  $L$ -Lipschitz w.r.t. the parameters  $v$

If we take  $m$  samples from  $p$  and  $q$  with

$$m \geq \frac{C n \Delta^2 \log(L L_\phi n / \epsilon)}{\epsilon^2}$$

$$\text{whp } |d_{\mathcal{F}, \phi}(\hat{p}, \hat{q}) - d_{\mathcal{F}, \phi}(p, q)| \leq \epsilon$$

Intuition: Both  $W_1$  and  $d_{F,\phi}$  were defined in terms of test functions, and for the latter we can bound the compl. using the # of parameters

Main Question: what can we say about the equilibrium?

Theorem [Arora et al., informal] If the generator has  $\tilde{\Omega}(n^2)$  parameters, then the generator can win

In particular, we will construct an approximate equilibrium where the discriminator can't do better than  $2\phi(\frac{1}{2})$

First we will need to introduce some tools from game theory

def: A zero sum game is defined by

- ① a set  $U$  of strategies for Alice
- ② a set  $V$  of strategies for Bob
- ③ A payoff function  $F(u,v)$  that for any  $u \in U, v \in V$  describes how much Alice wins / Bob loses

Notice the sum of payoffs is always equal to zero

Von Neumann's minimax theorem is a fundamental result in the area

Theorem [minimax for finite games]

If  $U$  and  $V$  are finite, there is a pair of distributions  $p$  and  $q$  on  $U$  and  $V$  respectively and a param.  $v$  called the game value s.t.

$$\textcircled{1} \quad \forall_{v \in V} \mathbb{E}_{u \sim p} [F(u, v)] \geq v = \text{payoff of } p \text{ and } q$$

$$\textcircled{2} \quad \forall_{u \in U} \mathbb{E}_{v \sim q} [F(u, v)] \leq v$$

Informally, Alice and Bob can each guarantee themselves  $\pm v$



Finally we call  $(p, a)$  an equilibrium

It is the solution to a convex/concave minimax problem, like we had for training GANs

$$\min_a \max_p \mathbb{E}_{\substack{u \sim p \\ v \sim q}} [F(u, v)] \quad (\text{or other way around})$$

Moreover if (1) and (2) hold only up to  $\pm \epsilon$ , we call it an  $\epsilon$ -approx. equilibrium

Note: There are versions of the minimax theorem that work with  $\infty$  strategy spaces, you need additional conditions

An important result of Lipton, Markakis and Mehta is:

Theorem 4 [Lipton et al.]: For any zero sum game with  $|U|=m, |V|=n$   
 $\exists \hat{p}, \hat{q}$  satisfying

①  $\hat{p}, \hat{q}$  is an  $\epsilon$ -approx. equilibrium

②  $\hat{p}$  has  $O\left(\frac{\log n}{\epsilon^2}\right)$  strategies in its support

$\hat{q}$  has  $O\left(\frac{\log m}{\epsilon^2}\right)$  "

The proof is by subsampling

Proof: Let  $(p, a)$  be an equilibrium.

Then let

$\hat{p}$  = empirical distributions of  
 $s = O\left(\frac{\log n}{\epsilon^2}\right)$  i.i.d draws from  $p$

and similarly for  $\hat{q}$ .

Then by standard concentration bds,  
no  $v \in V$  will do more than  $\epsilon$ -better on avg.  
playing against  $\hat{p}$  compared to  $p$

The same holds for any  $u \in U$  and  
 $\hat{q}$  and  $q$ .  $\square$

Now we will explain the main ingredients that go into analyzing GANs

First, you need the assumption

"The generator can approximate any point mass."

i.e. for any  $x$ ,  $\exists$  generator  $G_u$  with

$$\mathbb{P}_{y \sim D_{G_u}} [\|x - y\|] \leq \epsilon$$

Now the reasoning goes

① There is a mixed strategy for the generator that no discrim. can succeed against

This is true because you can represent  $D_{real}$

② There is a set of  $\tilde{O}(n)$  generators where the uniform distribution fools any discriminator

This relies on discretizing the strategy space of discriminators and subsampling mixed strategy as before

③ Can fold the generators into a larger deep net that uses the randomness in  $z$  to select uniformly from them

# Applications

Deep generative models can be used in downstream applications as more realistic models for real world inputs

Let's study compressed sensing

Setup: Unknown  $x \in \mathbb{R}^d$

We get linear measurements  $Ax = b$ ,  
where both  $A$  and  $b$  are known

How many rows (i.e. measurements) do we need to recover  $x$ ?

Claim: If  $x$  is arbitrary, then we need  $A$  to have full column rank

But what happens if  $x$  is structured (e.g. sparse)?

Suppose  $x$  has at most  $k$  nonzeros

Theorem [Donoho; Candes, Romberg, Tao]

For a random  $A$   <sup>$m \times d$</sup>  with <sup>standard normal</sup>

$$m \geq C k \log \frac{d}{k}$$

then whp can recover a  $k$ -sparse  $x$  exactly

In fact, the following algorithm

works:

$$\min \|z\|_1, \text{ s.t. } Az = b$$

and it is stable in the presence of noise / model misspecification

i.e. it'll approximately recover the  $k$  largest coordinates in  $x$

Implications: Can cutting radiation you're exposed to in an MRI by 90%

Food for Thought: Is sparsity a reasonable assumption for natural images?

It is a good starting point, but it turns out we can do much better!



In an influential work, Bora, Jalal, Price, Dimakis studied compressed sensing w/ generative models

Setup: Unknown  $x = G(z)$ , where  $G: \mathbb{R}^k \rightarrow \mathbb{R}^d$  is an  $L$ -layer network with ReLU activations

We observe  $Ax + n = b$

Main Question: How few measurements do we need to accurately recover  $x$ ?

They proved:

Theorem [Bora et al]: Suppose  $A$  is a random  $m \times d$  matrix with

standard normal entries. Further suppose

$$m \geq C L k \log d$$

Then whp the estimator

$$\hat{z} = \underset{z}{\operatorname{argmin}} \|b - AG(z)\|_2$$

satisfies  $\|G(\hat{z}) - x\|_2 \leq 6 \min_{z^*} \|G(z^*) - x\|_2 + 3 \|u\|_2$

Note:  $x$  can be anything — need not be a valid output of  $G$ . Thus we can approximate the best fit from few noisy linear measurements

Many of the usual notions from vanilla compressed sensing carry over with some twists

# Vanilla Compressed Sensing

Main Question: why can't there be two sparse solutions  $x$  and  $x'$  that both fit the data?

The key is that random matrices of the right dimensions satisfy the following conditions w.h.p

def: we say that  $A$  satisfies the  $(2k, \delta)$ -restricted isometry property if for all  $z$  with at most  $2k$  non zeros

$$(1-\delta) \|z\|_2^2 \leq \|Az\|_2^2 \leq (1+\delta) \|z\|_2^2$$

So if we have two sparse candidate solutions  $x \neq x'$ , they can't both fit our observations (even approximately)

For the sake of contradiction, suppose they do, i.e.

$$\|Ax - Ax'\|_2 \ll 1 - \delta$$

Then we can choose  $z = x - x'$  and

$$\|Az\|_2 \text{ is small}$$

which would violate the  $(2k, \delta)$ -RIP

Fact 2: A random  $m \times d$  matrix  $A$  with constant  $\times$  standard normal entries

$$m \geq ck \log \frac{d}{k} \Rightarrow (2k, \frac{1}{3})\text{-RIP whp}$$

Taking a step back:

RIP  
condition



compatibility btwn  
models for  $x$ 's and  
measurements

It turns out for more complex models,  
we can tweak the notion of RIP/REC

def. A matrix  $A$  satisfies the set  
restricted eigenvalue condition for a  
set  $S \subseteq \mathbb{R}^d$  with parameters  $\delta$  and  $\delta$   
if  $\forall x, x' \in S$  we have

$$\|A(x-x')\|_2 \geq \delta \|x-x'\|_2 - \delta$$

They set  $S = G(B^k(R))$   
ball of radius  $R$  in  $k$ -dimensions

Now the main ideas are

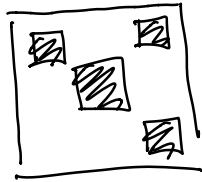
- ① A random matrix of appropriate  
size satisfies S-REC

This is the analogue of Fact 2

(2) S-REC guarantees that  $\hat{z}$  achieves the desired accuracy

But how do you find  $\hat{z}$ ? SGD

Applications: In image inpainting, you are missing many of the pixels



You can model this as observing

$$A \underbrace{G(z)}_{\text{natural image}}$$

where the rows of  $A$  correspond to observed pixels

By finding  $\hat{z}$ , can compute

$$\hat{x} = \text{inpainting} = G(\hat{z})$$

This leads to state-of-the-art results and some amazing pictures

Moral: Deep generative models can be a powerful replacement for simpler assumptions about the structure of realistic data