

## Problems w/ (Vanilla) GANs

- Convergence: min-max unstable

vanishing gradients

can't train D to completion

- if D too powerful, can overwhelm gen (100% separable, no useful gradients for G)

- if D too weak, random guessing also no gradients for G

GAN never converges, good performance quickly destabilized

- Model selection

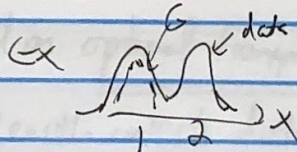
- D & G losses not correlated w/ quality

- in industry, / natural images often rely on "by eye" test!

- hard to know when to stop training

- Mode collapse

G & D stuck in vicious cycle



G learns to produce 1's

D fooled for a while

then realizes 2 are all real, guess on 2's

→ G will switch to making 2's

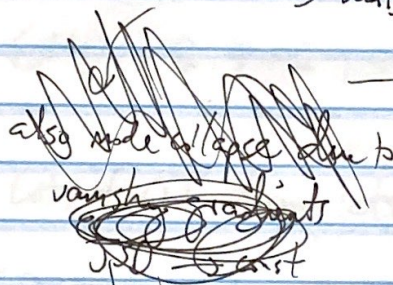
# Beyond vanilla GANs

- Many issues w/ GANs stem from loss fn

FD trained to optimality  $\rightarrow \min_G \text{JSD}(p_{data}, p_G)$

$\int p_{data} \leftrightarrow \int p_G \quad \text{JSD} = 1$

JSD = 1 whenever  $p_{data}$  &  $p_G$  disjoint  
 bounded from above  
 saturates  
 $\rightarrow$  Vanishing gradients!



$\rightarrow$  when G is very bad  $\rightarrow$  no ability to improve!

So can't train D to optimality for vanilla GAN

also causes mode collapse, stems from suboptimal D.

- Want: better GAN Loss, unbounded from above  
 informative even for widely separated distns

$\rightarrow$  one SOTA approach: Wasserstein GANs (1701.07875)

Arjovsky, Chintala, Bottou

Based on optimal transport theory

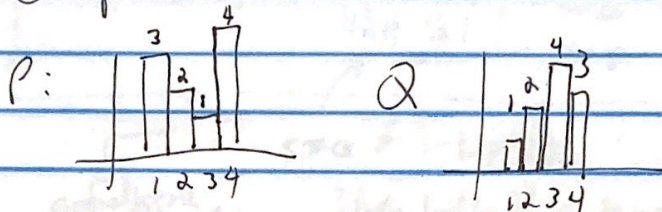
dist to distns  
 $P, Q$

"earth mover's distance" aka "Wasserstein distance"

"work it takes to transform P into Q"

$\hookrightarrow$  right property: more separated P & Q, more work

Example:



- transport 2 units  $1 \rightarrow 2$  so 1 matches
- transport 2 units  $2 \rightarrow 3$  so 2 matches
- transport 1 unit  $4 \rightarrow 3$  so 3 & 4 match ✓

EMD:  $2 \cdot 1 + 2 \cdot 1 + 1 \cdot 1 = 5$ .

↳ Lots of OT theory later...

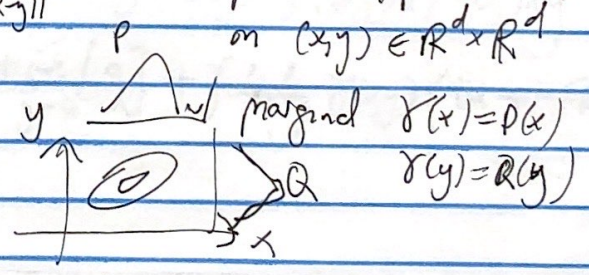
like MSE  
in prob  
distrib'n  
space

$$EMD(P, Q) = W(P, Q) = \min_{\gamma \in \Pi(P, Q)} \mathbb{E}_{(x, y) \sim \gamma} [\|x - y\|]$$

"Wasserstein-1"  
 Wasserstein-k is  $\mathbb{E}[\|x - y\|^k]$   
 space of all prob distrib'n  $\Pi(x, y)$  on  $(x, y) \in \mathbb{R}^d \times \mathbb{R}^d$

expectation value  
over all distrib'n  $\gamma$   
of distance

$W_1 = 0 \Leftrightarrow P = Q$   
unbounded for  $P \neq Q$   
 $\uparrow \Leftrightarrow \downarrow$  JS  $\rightarrow 1$   
 $W_1 \rightarrow \infty$



$W_1$  totally intractable...

→ Brilliant idea: use Kantorovich-Rubinstein duality

Amazing formula

K-R duality

$$W_1(P, Q) = \max_{\|h\|_L \leq 1} \left[ \mathbb{E}_{x \sim p} [h(x)] - \mathbb{E}_{y \sim q} [h(y)] \right]$$

space of 1-Lipschitz

actually true for K-Lipschitz  $\|h\|_L \leq K$

$$|h(x_1) - h(x_2)| \leq |x_1 - x_2| \quad \forall x_1, x_2$$

$$\Downarrow$$

$$|h'(x)| \leq 1 \text{ infinitesimal } \int \text{by mean value thm}$$

rough idea of pf: <sup>comes from</sup> ~~use~~ Lagrange multipliers

$$\begin{aligned} \mathbb{E}_{x \sim p, y \sim q} (|x - y|) &= \int dx dy \delta(x, y) (|x - y|) \\ &\quad + \int (p(x) - \int \delta(x, y) dy) f(x) dx \\ &\quad + \int (q(y) - \int \delta(x, y) dx) g(y) dy \\ &= \mathbb{E}_{x \sim p} [f] + \mathbb{E}_{y \sim q} [g] + \int dx dy \delta(x, y) (|x - y| - f(x) - g(y)) \end{aligned}$$

w/ K-R duality,  $W_1$  becomes tractable!

Approx  $h$  w/ NN  $\rightarrow$  "critic"  
analog of discriminator

$\max_{\|h\|_L \leq 1} (\dots) \rightarrow$  training the critic!