# Wasserstein GAN Martin Arjovsky, Soumith Chintala, and Leon Bottou

Presented by: Aditya Kumar Akash and Peyman Morteza

# Part I: Motivation and theoretical aspects of WGAN

# Unsupervised Learning Problem

- A sample data points from a hidden support is given.
- How to develop an algorithm to "learn" the hidden structure?



# Unsupervised Learning Problem

- structure has a complicated support?
- Why to study complexified version of the problem?



Image Data

• What if the underlying data points are in high dimension and the underlying



Music Data

# Problem formulation

- Let X be a compact metric space and let Σ denotes the Borel sets of X and P denotes the set of all probability measures on (X, Σ).
- Let  $p_r \in P$  be a hidden probability measure that generates the real data.
- Let  $X_p^n = \{x_1^p, \dots, x_n^p\}$  be a given i.i.d sample generated by  $p_r$ .



# Problem formulation

•  $p_{\theta}, \theta \in \Theta$ : a parametrized family of probability distributions.

• 
$$\theta^* := argmin_{\theta \in \Theta} (d(p_r, p_{\theta}))$$

• What is a good candidate for *d*?



## Looking for "metric" *d* such that $d(p_{\theta}, p_r) \rightarrow 0 \text{ as } \theta \rightarrow 0$

 $p_r$ : measure on  $[0,1]^2$  that is uniformly distributed on  $\{(0,y) \in [0,1]^2 | 0 \le y \le 1\}$ 





# "Distance" functions on P

$$KL(p \mid \mid q) := \int_{\mathcal{X}} \log(\frac{f_p(x)}{f_q(x)}) f_p(x) dx$$

 $JS(p,q) := KL(p,\frac{p+q}{2}) + KL(q,\frac{p+q}{2})$ 

 $\delta(p,q) := \sup_{K \in \Sigma} |p(K) - q(K)|$ 

Kullback–Leibler divergence: Shows up as asymptotic of MLE

Jensen–Shannon divergence: Symmetrized version of KL

Total variation between two measures

Looking for a "metric", d, such that  $d(p_{\theta}, p_r) \to 0$  as  $\theta \to 0$ If  $\theta > 0 \implies p_r$  and  $p_{\theta}$  have disjoint support  $\implies$   $KL(p_r, p_{\theta}) = \infty \implies KL$  is not even defined

in this particular example.

 $p_r$ : measure on  $[0,1]^2$  that is uniformly distributed on  $\{(0,y) \in [0,1]^2 | 0 \le y \le 1\}$ 





# **JS Divergence**

Looking for a "metric", d, such that  $d(p_{\theta}, p_r) \to 0$  as  $\theta \to 0$ If  $\theta > 0 \implies \frac{p_r + p_{\theta}}{2}$  has mass  $\frac{1}{2}$  on  $y = 0, y = \theta$  $\implies JS(p_r, p_{\theta}) = \frac{1}{2}\log(2) + \frac{1}{2}\log(2) = \log(2)$  $\implies$  JS is constant and does not change by  $\theta$ .

 $p_r$ : measure on  $[0,1]^2$  that is uniformly distributed on  $\{(0,y) \in [0,1]^2 | 0 \le y \le 1\}$ 





# **Total Variation**

Looking for a "metric", d, such that  $d(p_{\theta}, p_r) \to 0$  as  $\theta \to 0$ 

If  $\theta > 0 \implies \delta(p_r, p_{\theta}) = 1$  is constant and does not change by  $\theta$ .

If  $\theta = 0 \implies \delta(p_r, p_{\theta}) = 0.$ 

 $p_r$ : measure on  $[0,1]^2$  that is uniformly distributed on  $\{(0,y) \in [0,1]^2 | 0 \le y \le 1\}$ 





- Transport each blue ball and match it to a yellow ball.
- c(x, y) = |x y| cost of moving a ball located at position x to position y. • Find a transportation plan with minimal cost.



### Optimal Transportation Plan







# **Transportation Plan**

A transportation plan can be viewed as a probability measure on  $\mathcal{X} \times \mathcal{X}$  with specified marginals.

 $\begin{array}{ccc}
\vdots \\
\mu(x,y) = \begin{cases} \frac{1}{2} & (x,y) = (1,2) \\
\frac{1}{3} & (x,y) = (3,4), \\
\frac{1}{6} & (x,y) = (5,4) \\
\end{array}$ 

 $c(x, y)d\mu(x, y) = 1$  $JX \times X$ 

Number of balls that is transported from position 3 to position 4.





# Wasserstein Distance

# $\Pi(p,q) := \{ \mu \in Prob(\mathcal{X} \times \mathcal{X}) | (\pi_1)_*(\mu) = p, (\pi_2)_*(\mu) = q \}.$ $c: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ $c(x, y) = d_{\mathcal{X}}(x, y)$ $W(p,q) := \inf_{\mu \in \Pi(p,q)} \int_{\mathcal{X} \times \mathcal{X}} c(x,y) d\mu$

The set of all transportation plans between p and q

Wasserstein distance is the infimum cost among all transportation plans



# Wasserstein Distance

## Looking for a "metric", *d*, such that $d(p_{\theta}, p_r) \rightarrow 0$ as $\theta \rightarrow 0$

## $W(p_r, p_{\theta}) = \theta \to 0.$

 $p_r$ : measure on  $[0,1]^2$  that is uniformly distributed on  $\{(0,y) \in [0,1]^2 | 0 \le y \le 1\}$ 

 $p_{\theta}: \text{ measure on } [0,1]^2 \text{ that is uniformly}$ distributed on  $\{(\theta, y) \in [0,1]^2 \mid 0 \le y \le 1, \theta \in (0,1)\}$ 



## **General case formulation**



# Formal Theorem

### Theorem.

- If g is continuous in  $\theta$  then  $\phi$  is continuous.

then  $\phi$  is differentiable almost everywhere.

Theorem 1 in the paper.

everywhere differentiable.

• If g is locally Lipschitz and the local constants  $L(\theta, z)$  satisfies,

$$\int_{\mathcal{Z}} L(\theta, z) d\mu(z) < \infty$$

**Corollary.** If  $g_{\theta}$  is a neural net and  $\int_{\mathcal{Z}} ||z|| d\mu(z) < \infty$  then  $\phi(\theta)$  is continuous and almost

Corollary 1 in the paper



# Optimizing $\phi$

## $\phi(\theta) := W(p_{\theta}, p_{r})$

# $W(p_{\theta}, p_{r}) := \inf_{\mu \in \Pi(p_{\theta}, p_{r})} \int_{\mathcal{X} \times \mathcal{X}} c(x, y) d\mu$

 $\Pi(p_r, p_\theta) := \{ \mu \in Prob(\mathcal{X} \times \mathcal{X}) \, | \, (\pi_1)_*(\mu) = p_r, (\pi_2)_*(\mu) = p_\theta \} \, .$ 



Not obvious from this representation how to compute  $\phi$ in general or differentiate it.

## Kantorovich-Rubinstein duality

$$\phi(\theta) := W(p_r, p_\theta) = \sup_{\psi \in Lip_1} \left( \int_{\mathcal{X}} \left( \int$$

Theorem 3 in WGAN paper

$$f^* = \operatorname{argmax}_{\psi \in Lip_1} \left( \int_{\mathcal{X}} \psi(x) dp_r(x) \right)$$
$$\frac{\partial \phi}{\partial \theta} = -\int_{\mathcal{Z}} \frac{\partial (f^* \circ g_\theta)(z)}{\partial \theta} d\mu(z)$$

### No intractable inf anymore!

 $(x)dp_r(x) - \int_{\infty} \psi(x)dp_{\theta}(x)$ 

When both sides are defined



 $\mu$ : your favorite measure on  $\mathcal{X}$  (e.g. gaussian, uniform, etc)

 $g: \mathcal{X} \times \Theta \to \mathcal{X}$ 

 $\mathcal{I}$  .

 $p_{\theta} := (g_{\theta})_* \mu$ 

 $\phi(\theta) := W(p_{\theta}, p_{r})$ 

 $x^i \sim (\mathcal{X}, p_r)$ 

# WGAN Algorithm



 $x_i \sim (\mathcal{X}, p_r)$   $z_i \sim (\mathcal{Z}, \mu)$ 

 $\phi(\theta) := W(p$  $g_{\theta}: \mathcal{X} \to \mathcal{X}, \theta \in \Theta$ 

A neural net that generates samples

Want to minimize with gradient descent

Needs to be compact  $f_w: \mathcal{X} \to \mathbb{R}, w \in \mathcal{W}$  $f_w(x)$ 

Second neural net for estimating  $f^*$ 

Want to maximize with gradient ascent

# WGANAlgorithm

$$\mathcal{P}_{r}, \mathcal{P}_{\theta}) \qquad \frac{\partial \phi}{\partial \theta} = -\int_{\mathscr{Z}} \frac{\partial (f^{*} \circ g_{\theta})(z)}{\partial \theta} d\mu(z) \quad \frac{\partial}{\partial \theta} (\frac{1}{n} \sum_{i=1}^{n} f^{*} \circ g_{\theta}) d\mu(z)$$
Want to estimate Need to estimate  $f^{*}$ 

$$f_{\mathcal{X}}(x) - \int_{\mathcal{X}} f_{w}(x) dp_{\theta}(x) \Big)$$

$$\frac{\partial}{\partial w} \Big( \sum_{i=1}^{n} \frac{1}{n} f_w(x^i) - \frac{1}{n} \sum_{i=1}^{n} f_w(g_e) \Big)$$

We can compute gradient of the estimation with respect

to *W* 





Algorithm 1 WGAN, our proposed algorithm. All experiments in the paper used the default values  $\alpha = 0.00005$ , c = 0.01, m = 64,  $n_{\text{critic}} = 5$ . **Require:** :  $\alpha$ , the learning rate. c, the clipping parameter. m, the batch size.  $n_{\rm critic}$ , the number of iterations of the critic per generator iteration. **Require:** :  $w_0$ , initial critic parameters.  $\theta_0$ , initial generator's parameters. 1: while  $\theta$  has not converged do for  $t = 0, ..., n_{\text{critic}}$  do 2: Sample  $\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$  a batch from the real data. 3: Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples. 4:  $g_w \leftarrow \nabla_w \left[ \frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)})) \right]$ 5: $w \leftarrow w + \alpha \cdot \text{RMSProp}(w, g_w)$ 6:  $w \leftarrow \operatorname{clip}(w, -c, c)$ 7:end for 8: Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples. 9:  $g_{\theta} \leftarrow -\nabla_{\theta} \frac{1}{m} \sum_{i=1}^{m} f_w(g_{\theta}(z^{(i)}))$ 10:  $\theta \leftarrow \theta - \alpha \cdot \text{RMSProp}(\theta, q_{\theta})$ 11:12: end while

# Part II: Training and empirical aspects of WGAN

# **Training WGAN**

## **A better GAN loss function**

- EM distance is continuous and differentiable a.e
- Reliable gradients with better trained critic



## No mode collapse for WGAN

Mode collapse for Standard GAN



![](_page_26_Picture_4.jpeg)

# Training WGAN

## Explanations for mode collapse?

- Standard GAN
  - For a fixed discriminator, generator uses few modes to fool it
  - Training critic til optimality provides less meaningful gradients
- WGAN
  - For critic trained till optimality, we get more reliable gradients

# **Empirical Results**

![](_page_28_Picture_1.jpeg)

## Datasets & Baselines

- Image generation Datasets
  - LSUN-Bedroom dataset

- Baselines
  - DCGAN GAN with conv architectures

![](_page_29_Picture_5.jpeg)

# Standard GAN Training

![](_page_30_Figure_1.jpeg)

### Objective

![](_page_30_Picture_3.jpeg)

![](_page_30_Figure_4.jpeg)

## $\min_{C} \max_{D} V(D, G) = \mathbb{E}_{x \sim p_{data}(x)}[\log(D(x))] + \mathbb{E}_{z \sim p_{z}(z)}[\log(1 - D(G(z)))]$

## **DCGAN Generator**

- Replace pooling layers with convolutions
- Use batchnorm
- Remove FC hidden layers
- Use ReLU for G
- Use Leaky ReLU for D

![](_page_31_Picture_6.jpeg)

# WGAN as Meaningful loss

![](_page_32_Figure_2.jpeg)

For DCGAN generator sample quality better,

Increasing error (JS Estimate) for MLP generator (FC Network) But no correlation to JS estimates

For MLP gen and disc, Curve fluctuates regardless of quality

# WGAN as Meaningful loss

### Wasserstein estimates with WGAN training

![](_page_33_Figure_2.jpeg)

![](_page_34_Picture_0.jpeg)

### **Negative Result**

- Training unstable with
  - Momentum based optimizer eg ADAM
  - High learning rates
- Since critic loss is non-stationary
  - momentum a potential cause
  - Works with RMSProp -

# WGAN Training

$$E[g^{2}]_{t} = \beta E[g^{2}]_{t-1} + (1-\beta)(\frac{\delta C}{\delta w})^{2}$$
$$w_{t} = w_{t-1} + (\frac{\eta}{\sqrt{E[g^{2}]_{t}}} \frac{\delta C}{\delta w}$$

Adaptive LR

# Improved Stability - WGAN

## WGAN vs Standard GAN

WGAN training

![](_page_35_Picture_3.jpeg)

![](_page_35_Picture_4.jpeg)

Standard GAN training

### High quality samples for both

# Improved Stability - WGAN

## WGAN vs Standard GAN

### WGAN training - Able to produce samples

Standard GAN training - fails

![](_page_36_Picture_4.jpeg)

![](_page_36_Picture_5.jpeg)

DCGAN Generator without batch normalization and constant number of filter each layer

# Improved Stability - WGAN

### WGAN vs Standard GAN

WGAN training - Able to produce good samples

Standard GAN training - Lower quality images

![](_page_37_Picture_4.jpeg)

![](_page_37_Picture_5.jpeg)

Notice the mode collapse for standard GAN samples

### MLP Generator

# Related Works

![](_page_38_Picture_1.jpeg)

## Works with Integral Probability Metrics (IPM)

- IPM with class  $\mathscr{F}$ :  $d_{\mathscr{F}}(\mathbb{P}_r, \mathbb{P}_{\theta}) = \sup \mathbb{E}_{x \sim \mathbb{P}_r}[f(x)] \mathbb{E}_{x \sim \mathbb{P}_{\theta}}[f(x)]$ f∈ℱ
- $\mathcal{F}$ : 1-Lipschitz gives  $W(\mathbb{P}_r, \mathbb{P}_{\rho})$
- $\mathcal{F}: \{f: f(x) \in [0,m] \forall x\}$  for Energy-based GANs (EBGANs)
  - Same behavior as Total variation distance

# • $\mathcal{F}: \{f: f(x) \in [-1,1] \forall x\}$ gives $d_{\mathcal{F}}(\mathbb{P}_r, \mathbb{P}_{\theta}) = \delta(\mathbb{P}_r, \mathbb{P}_{\theta})$ (Total variation distance)

# Works with Integral Probability Metrics (IPM)

- $\mathcal{F} = \{ f \sim \mathcal{H} : ||f||_{\mathcal{H}} \leq 1 \}$ , for some RKHS  $\mathcal{H}$  associated with kernel  $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ gives Maximum Mean Discrepancy (MMD)
  - MMD is a proper metric & no need to maximize the inner supremum
  - Disadvantage: quadratic computational cost
  - Linear approximations have worse sample complexity
- Generative Moment Matching Networks (GMMs): Generative networks using MMD
  - Quadratic cost as a function of samples
  - Vanishing gradients for low-bandwidth kernels

# Follow-up Works

![](_page_41_Picture_1.jpeg)

# WGAN with Gradient Penalty

- Difficulties with weight constraints:

Weight clipping critic fail to capture higher moments

Critic with Gradient Penalty

![](_page_42_Figure_5.jpeg)

![](_page_42_Figure_6.jpeg)

![](_page_42_Figure_7.jpeg)

Value surface of WGAN critics

## • Optimization difficulties - Resulting critic can have a pathological value surface

### 25 Gaussians

![](_page_42_Figure_11.jpeg)

![](_page_42_Picture_12.jpeg)

Swiss Roll

![](_page_42_Picture_13.jpeg)

![](_page_42_Picture_14.jpeg)

# WGAN with Gradient Penalty

$$L = \underbrace{\mathbb{E}_{\tilde{\boldsymbol{x}} \sim \mathbb{P}_{g}} \left[ D(\tilde{\boldsymbol{x}}) \right] - \mathbb{E}_{\boldsymbol{x} \sim \mathbb{P}_{r}} \left[ D(\boldsymbol{x}) \right]}_{\text{Original critic loss}}$$

• Stable training

• A differentiable function is 1-Lipschitz iff it has gradients with norm at most 1 everywhere.

 $+ \lambda \mathop{\mathbb{E}}_{\hat{\boldsymbol{x}} \sim \mathbb{P}_{\hat{\boldsymbol{x}}}} \left[ (\|\nabla_{\hat{\boldsymbol{x}}} D(\hat{\boldsymbol{x}})\|_2 - 1)^2 \right].$ Our gradient penalty Soft constraint

# **Sliced WGAN**

- Wasserstein computation in dual involves solving saddle point problem

Sliced Wasserstein:

![](_page_44_Picture_6.jpeg)

• Another variant of Wasserstein uses random projections and is a single minimization

# Summary

- Introduced Wasserstein distances for GANs
  - like KL, JS
- Define a form of GAN using duality : Wasserstein GAN
  - Meaningful loss metric for distributions
  - Stable training of GAN
  - Gets rid of issues like mode collapse
- Experiments show stable learning of image generators

### • Theoretical analysis of EM distance behavior in comparison to popular divergences

# Thank You!

![](_page_46_Picture_1.jpeg)

# Questions?

![](_page_47_Picture_1.jpeg)

# Quiz Problems

- KL divergence
- JS divergence
- Total variation distance(TV)
- Wasserstein distance (Earth-Mover)

According to the paper, which of the following gives a more sensible cost function for learning probability distribution that are supported in low dimensional submanifolds?

![](_page_48_Picture_7.jpeg)

# Quiz Problems

- In WGAN, we can train the critic till optimality.
- WGAN solves the mode collapse issue in GAN.
- WGAN training is not stable when one uses momentum based optimizer or high learning rate.
- All of the above.

![](_page_49_Picture_6.jpeg)

Which of the following statements is true:

![](_page_50_Picture_0.jpeg)

- Experiments illustrate that Wasserstein estimates correlates well with the quality of the generated samples.
- WGAN training is reported to be stable with Adam optimizer.
- It is observed that WGANs are much more stable than GANs when one varies the generator architecture.
- JS estimates follow the same correlation with image quality as the Wasserstein estimates.

# Quiz Problems

Which of the following found to be true in empirical evaluation of WGANs?

# References

- The slides have materials from following works:
- Wasserstein GANs <u>https://arxiv.org/abs/1701.07875</u>
- Generative Adversarial Networks https://arxiv.org/abs/1406.2661
- lacksquare
- Unrolled Generative Adversarial Networks <u>https://arxiv.org/abs/1611.02163</u>
- Generative Modelling using the Sliced Wasserstein Distance <u>https://arxiv.org/abs/1803.11188</u>
- Improved Training of Wasserstein GANs <u>https://arxiv.org/abs/1704.00028</u>
- LSUN Dataset: <u>https://www.yf.io/p/lsun</u>
- Image Sources at slide 4: https://www.pinterest.com/pin/382383824583177634/ and https://www.pinterest.com/pin/  $\bullet$ 161777811601964865/

Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks - <u>https://arxiv.org/abs/1511.06434</u>

Understanding RMSprop - https://towardsdatascience.com/understanding-rmsprop-faster-neural-network-learning-62e116fcf29a