

Wasserstein Distance in Machine Learning

Huidong (Leon) Liang

Discipline of Business Analytics
The University of Sydney Business School

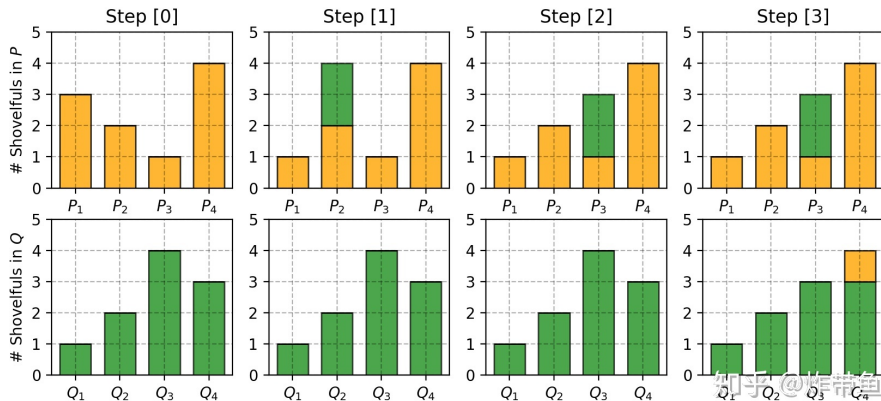
November 2021

- 1 Wasserstein Distance
- 2 Wasserstein GAN
- 3 Wasserstein Auto-Encoders

Wasserstein Distance

Optimal Transport

Earth Mover Distance



知乎@炸带鱼

Wasserstein Distance and Kantorovich-Rubinstein Duality

Definition

- 1-Wasserstein distance between two distributions \mathbb{P}_r and \mathbb{P}_g :

$$W_1(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \mathcal{P}(r \sim \mathbb{P}_r, z \sim \mathbb{P}_g)} \mathbb{E}[\|r - z\|_2] \quad (1)$$

where $\mathcal{P}(r \sim \mathbb{P}_r, z \sim \mathbb{P}_g)$ is the set of all joint distributions $\gamma(r, g)$ with marginals \mathbb{P}_r and \mathbb{P}_g respectively.

- Kantorovich-Rubinstein duality:

$$W_1(\mathbb{P}_r, \mathbb{P}_g) = \sup_{\|f\|_L \leq 1} \{ \mathbb{E}_{r \sim \mathbb{P}_r}[f(r)] - \mathbb{E}_{z \sim \mathbb{P}_g}[f(z)] \} \quad (2)$$

where f is any continuous function that satisfies 1-Lipschitz continuity.

Wasserstein Distance, KL divergence and JS divergence

Comparison

- Kullback-Leibler Divergence

$$\mathbf{KL}(\mathbb{P}_r || \mathbb{P}_g) = \int \log\left(\frac{\mathbb{P}_r(x)}{\mathbb{P}_g(x)}\right) \mathbb{P}_r(x) dx \quad (3)$$

- Jensen-Shannon Divergence

$$\mathbf{JS}(\mathbb{P}_r || \mathbb{P}_g) = \frac{1}{2} \mathbf{KL}(\mathbb{P}_r || \frac{\mathbb{P}_g + \mathbb{P}_r}{2}) + \frac{1}{2} \mathbf{KL}(\mathbb{P}_g || \frac{\mathbb{P}_g + \mathbb{P}_r}{2}) \quad (4)$$

- 1-Wasserstein Distance

$$W_1(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \mathcal{P}(r \sim \mathbb{P}_r, z \sim \mathbb{P}_g)} \mathbb{E}[\|r - z\|_2] \quad (5)$$

Wasserstein Distance, KL divergence and JS divergence

Comparison

- Let $z \sim U[0, 1]$ be a random variable from an 1-dimensional uniform distribution. Let \mathbb{P}_0 be a distribution of $(0, z) \in \mathbb{R}^2$ (i.e. 0 on the x-axis and z on the y-axis), which is uniform on a vertical line that passes through the origin. We now define $\mathbb{P}_\theta(z)$ of (θ, z) with one single parameter θ [1]. It is not hard to find the following:

- $W(\mathbb{P}_0, \mathbb{P}_\theta) = |\theta|,$

- $JS(\mathbb{P}_0, \mathbb{P}_\theta) = \begin{cases} \log 2 & \text{if } \theta \neq 0, \\ 0 & \text{if } \theta = 0, \end{cases}$

- $KL(\mathbb{P}_\theta \| \mathbb{P}_0) = KL(\mathbb{P}_0 \| \mathbb{P}_\theta) = \begin{cases} +\infty & \text{if } \theta \neq 0 \\ 0 & \text{if } \theta = 0 \end{cases}$

Wasserstein Distance, KL divergence and JS divergence

Comparison

- Although this simple example features distributions with disjoint supports, the same conclusion holds when the supports have a non empty intersection contained in a set of measure zero. This happens to be the case when two low dimensional manifolds intersect in general position [1].

Wasserstein GAN

Related Work: GAN [2]

Loss Function and Algorithm

- Objective Function

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

- Algorithm

for number of training iterations **do**

for k steps **do**

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

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

end for

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

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

end for

Related Work: GAN

The Optimal Discriminator

- the Optimal Discriminator

$$-P_r(x) \log D(x) - P_g(x) \log[1 - D(x)]$$

$$-\frac{P_r(x)}{D(x)} + \frac{P_g(x)}{1 - D(x)} = 0$$

$$D^*(x) = \frac{P_r(x)}{P_r(x) + P_g(x)}$$

- Substitute back to the objective function

$$\mathbb{E}_{x \sim P_r} \log \frac{P_r(x)}{\frac{1}{2}[P_r(x) + P_g(x)]} + \mathbb{E}_{x \sim P_g} \log \frac{P_g(x)}{\frac{1}{2}[P_r(x) + P_g(x)]} - 2 \log 2$$

$$2JS(P_r || P_g) - 2 \log 2$$

Wasserstein GAN

Loss Function

- Kantorovich-Rubinstein duality:

$$W_1(\mathbb{P}_r, \mathbb{P}_g) = \sup_{\|f\|_L \leq 1} \{ \mathbb{E}_{r \sim \mathbb{P}_r}[f(r)] - \mathbb{E}_{z \sim \mathbb{P}_g}[f(z)] \} \quad (6)$$

where f is any continuous function that satisfies 1-Lipschitz continuity.

- Loss Function of Wasserstein GAN [1]

$$\max_{w \in \mathcal{W}} \mathbb{E}_{x \sim \mathbb{P}_r}[f_w(x)] - \mathbb{E}_{z \sim p(z)}[f_w(g_\theta(z))]$$

w is clipped in to a compact space \mathcal{W} for Lipschitz constraint.

- gradient penalty for Lipschitz constraint: add $\lambda \cdot \mathbb{E}(\|\nabla f(x)\| - 1)^2$ to the objective.

Wasserstein GAN

- Wasserstein GAN Algorithm

```

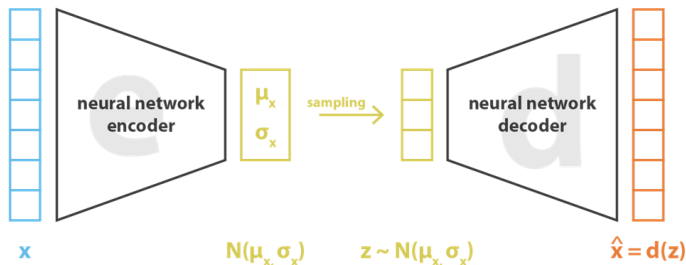
1: while  $\theta$  has not converged do
2:   for  $t = 0, \dots, n_{\text{critic}}$  do
3:     Sample  $\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$  a batch from the real data.
4:     Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of priors.
5:      $g_w \leftarrow \nabla_w [\frac{1}{m} \sum_{i=1}^m f_w(x^{(i)})$ 
         $\quad \quad \quad - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))]$ 
6:      $w \leftarrow w + \alpha \cdot \text{RMSProp}(w, g_w)$ 
7:      $w \leftarrow \text{clip}(w, -c, c)$ 
8:   end for
9:   Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
10:   $g_\theta \leftarrow -\nabla_\theta \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))$ 
11:   $\theta \leftarrow \theta - \alpha \cdot \text{RMSProp}(\theta, g_\theta)$ 
12: end while

```

Wasserstein Auto-Encoders

Related Work: VAE [3]

- Model Structure



- Loss function

$$\mathcal{L} = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log p(\mathbf{x}|\mathbf{z}) \right] - \mathbf{KL} \left(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}) \right) \quad (7)$$

Wasserstein Auto-Encoders

Intuition

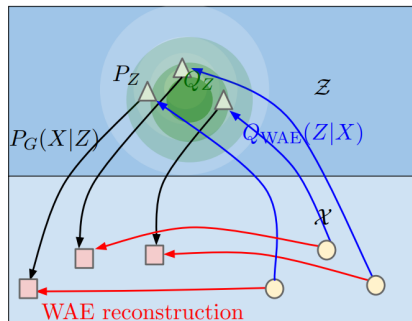
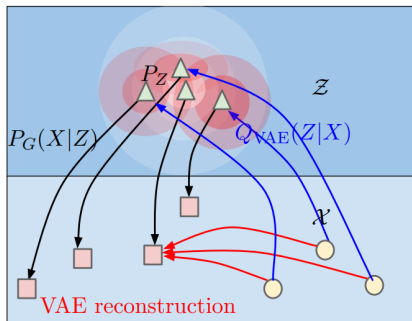


Figure from Wasserstein Auto-Encoders [4]

Wasserstein Auto-Encoders

Objective Function

- Objective Function

$$\inf_{\Gamma \in \mathcal{P}(X \sim P_X, Y \sim P_G)} \mathbb{E}_{(X,Y) \sim \Gamma} [c(X, Y)] = \inf_{Q: Q_Z = P_Z} \mathbb{E}_{P_X} \mathbb{E}_{Q(Z|X)} [c(X, G(Z))]$$

where Q_Z is the marginal distribution of Z when $X \sim P_X$ and $Z \sim Q(Z|X)$

$$p_G(x) := \int_Z p_G(x|z) p_Z(z) dz, \quad \forall x \in \mathcal{X}$$

- Relax the constraints on Q_Z by adding a penalty term

$$D_{\text{WAE}}(P_X, P_G) := \inf_{Q(Z|X) \in \mathcal{Q}} \mathbb{E}_{P_X} \mathbb{E}_{Q(Z|X)} [c(X, G(Z))] + \lambda \cdot \mathcal{D}_Z(Q_Z, P_Z)$$

Wasserstein Auto-Encoders

GAN-based Penalty (WAE-GAN) and MMD-based penalty (WAE-MMD)

- GAN-based Penalty (WAE-GAN)

$$\mathcal{D}_Z(Q_Z, P_Z) = D_{\text{JS}}(Q_Z, P_Z)$$

- Maximum Mean Discrepancy (MMD)-based Penalty (WAE-MMD)

$$\text{MMD}_k(P_Z, Q_Z) = \left\| \int_{\mathcal{Z}} k(z, \cdot) dP_Z(z) - \int_{\mathcal{Z}} k(z, \cdot) dQ_Z(z) \right\|_{\mathcal{H}_k}$$

The End

Reference



Martin Arjovsky, Soumith Chintala, and Léon Bottou.

Wasserstein generative adversarial networks.

In *International Conference on Machine Learning*, pages 214–223. PMLR, 2017.



Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio.

Generative adversarial nets.

Advances in Neural Information Processing Systems, 27, 2014.



Diederik P Kingma and Max Welling.

Auto-encoding variational Bayes.

In *International Conference on Learning Representations*, 2014.



Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schoelkopf.

Wasserstein auto-encoders.

In *International Conference on Learning Representations*, 2018.