

Probabilistic Deep Learning for Solving Inverse Problems

Uncertainty Quantification Summer School

Assad A Oberai

USC

August 4th 2023

Includes work with

- Deep Ray
- Javier Mungoitio Esandi
- Harisankar Ramaswamy
- Bryan Shadday
- Agnimitra Dasgupta

Outline

Introduction

Wasserstein Generative Adversarial Network

Conditional Wasserstein Generative Adversarial Network

Generative Diffusion Networks

Conditional Generative Diffusion Networks

Table of Contents

Introduction

Wasserstein Generative Adversarial Network

Conditional Wasserstein Generative Adversarial Network

Generative Diffusion Networks

Conditional Generative Diffusion Networks

Generative Problem

Given $\mathbf{x}^{(i)} \sim P_{\text{data}}(\mathbf{x})$ generate more samples.

Utility:

- Data augmentation.
- New designs.
- Reduced representation.

Conditional Generative Problem

Let \mathbf{X} denote a random vector to be inferred, and \mathbf{Y} denote the measured random vector. Given

1. The measurement $\mathbf{Y} = \hat{\mathbf{y}}$.
2. Some prior knowledge of \mathbf{X} encoded in $P_{\text{data}}(\mathbf{x})$.
3. The forward map $\mathbf{y} = \mathbf{f}(\mathbf{x}, \eta)$.

generate samples from $P_{\text{data}}(\mathbf{x}|\hat{\mathbf{y}})$.

Utility:

- Inverse problems with a known forward map.
- Data assimilation.

Can be reformulated

Conditional Generative Problem

Steps

1. Sample $\mathbf{x}^{(i)} \sim P_{\text{data}}(\mathbf{x})$.
2. Compute $\mathbf{y}^{(i)} = \mathbf{f}(\mathbf{x}^{(i)}, \eta)$.
3. Collect $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \sim P_{\text{data}}(\mathbf{x}, \mathbf{y})$.

Reduced problem: given

- The measurement $\mathbf{Y} = \hat{\mathbf{y}}$.
- $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \sim P_{\text{data}}(\mathbf{x}, \mathbf{y})$

generate samples from $P_{\text{data}}(\mathbf{x}|\hat{\mathbf{y}})$.

- Wasserstein Generative Adversarial Networks (WGAN).
- Conditional Wasserstein Generative Adversarial Networks (CWGAN).
- Generative Diffusion Networks.
- Conditional Generative Diffusion Networks.

Table of Contents

Introduction

Wasserstein Generative Adversarial Network

Conditional Wasserstein Generative Adversarial Network

Generative Diffusion Networks

Conditional Generative Diffusion Networks

Architecture

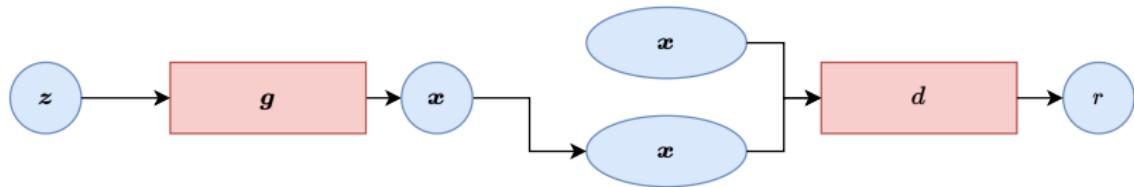


Figure: Schematic of a WGAN.

Key references:

- Villani, Cédric. Optimal transport: old and new. Vol. 338. Berlin: Springer, 2009.
- Arjovsky, Martin, Soumith Chintala, and Léon Bottou. "Wasserstein generative adversarial networks." International conference on machine learning. PMLR, 2017.
- Gulrajani, Ishaan, et al. "Improved training of wasserstein gans." Advances in neural information processing systems 30 (2017).

WGAN: Generator and Critic

Properties of the generator:

- $\mathbf{g} : \mathbb{R}^{N_z} \times \mathbb{R}^{N_\theta} \rightarrow \mathbb{R}^{N_x}$
- $\mathbf{z} \sim P_Z$ and $N_z \ll N_x$
- $\mathbf{x} \sim P_g \Rightarrow \mathbf{x} = \mathbf{g}(\mathbf{z}; \theta), \mathbf{z} \sim P_Z.$

Properties of the critic:

- $d : \mathbb{R}^{N_x} \times \mathbb{R}^{N_\phi} \rightarrow \mathbb{R}, \quad d : \mathbf{x} \mapsto r$
- $r = d(\mathbf{x}; \phi).$

“Loss” function:

$$L(\theta, \phi) = \mathbb{E}_{P_{\text{data}}} [d(\mathbf{x}; \phi)] - \mathbb{E}_{P_Z} [d(\mathbf{g}(\mathbf{z}; \theta); \phi)].$$

Optimality

Optimal values (min-max problem):

$$(\theta^*, \phi^*) = \operatorname{argmin}_{\theta} \operatorname{argmax}_{\phi, d \in \text{Lip}} L(\theta, \phi).$$

With infinite data, infinite critic capacity achieving the maximum implies,

$$\theta^* = \operatorname{argmin}_{\theta} W_1(P_{\text{data}}, P_g).$$

If this minima is achieved $W_1(P_{\text{data}}, P_{g^*}) = 0$. This implies weak convergence of the measures. Which in turn means,

$$\mathbb{E}_{P_{\text{data}}} [I(x)] = \mathbb{E}_{P_Z} [I(g(z; \theta^*))], \forall I \in C_b(\Omega_X).$$

Implementation of WGAN

- In practice, the “loss” term is approximated as:

$$L(\theta, \phi) = \frac{1}{N} \left(\sum_{i=1}^N d(\mathbf{x}^{(i)}; \phi) - \sum_{i=1}^N d(\mathbf{g}(\mathbf{z}^{(i)}; \theta); \phi) \right) + \lambda P(\phi).$$

$$P(\phi) = \frac{1}{N} \sum_{i=1}^N (\|\nabla d(\mathbf{h}^{(i)}; \phi)\|_2 - 1)^2$$

$$\mathbf{h}^{(i)} = \delta^{(i)} \mathbf{x}^{(i)} + (1 - \delta^{(i)}) \mathbf{g}(\mathbf{z}^{(i)}; \theta), \quad \delta^{(i)} \sim U(0, 1).$$

- The min-max problem is hard to solve. Use stochastic gradient descent-ascent. Five steps of ascent followed by one descent.
- Hard to track convergence. Look for stability of the “loss” function and track some statistic of the generated sample.

Implementation of WGAN

- If x is a vector, the generator and the critic are MLPs.
- If x is an image (discrete field), the generator is a U-Net with transpose convolutions. The critic is a CNN - like a classifier.

Table of Contents

Introduction

Wasserstein Generative Adversarial Network

Conditional Wasserstein Generative Adversarial Network

Generative Diffusion Networks

Conditional Generative Diffusion Networks

Architecture

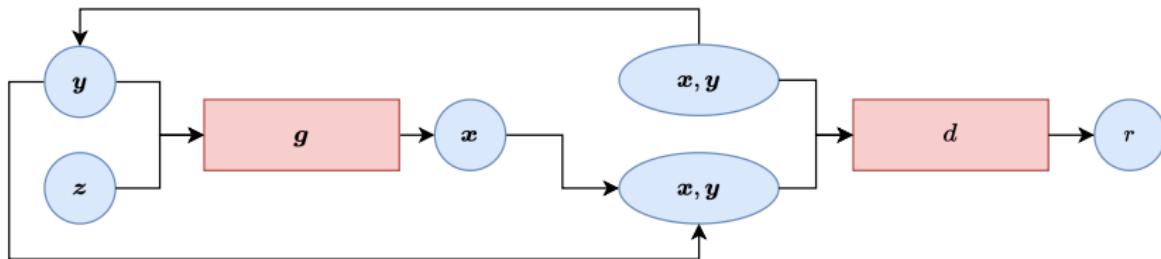


Figure: Schematic of a conditional WGAN.

Key references:

- Adler, Jonas, and Ozan Öktem. "Deep bayesian inversion." arXiv preprint arXiv:1811.05910 (2018).
- Ray, Deep, et al. "The efficacy and generalizability of conditional GANs for posterior inference in physics-based inverse problems." arXiv preprint arXiv:2202.07773 (2022).
- Ray, Deep, et al. "Solution of physics-based inverse problems using conditional generative adversarial networks with full gradient penalty." arXiv preprint arXiv:2306.04895 (2023).

CWGAN: Generator and Critic

Properties of the generator:

- $\mathbf{g} : \mathbb{R}^{N_z} \times \mathbb{R}^{N_y} \times \mathbb{R}^{N_\theta} \rightarrow \mathbb{R}^{N_x}$
- $\mathbf{z} \sim P_Z$ and $N_z \ll N_x$
- $\mathbf{x} \sim P_g(\mathbf{x}|\mathbf{y}) \Rightarrow \mathbf{x} = \mathbf{g}(\mathbf{z}, \mathbf{y}; \theta), \mathbf{z} \sim P_Z.$

Properties of the critic:

- $d : \mathbb{R}^{N_x} \times \mathbb{R}^{N_y} \times \mathbb{R}^{N_\phi} \rightarrow \mathbb{R}$
- $r = d(\mathbf{x}, \mathbf{y}; \phi).$

“Loss” function:

$$L(\theta, \phi) = \mathbb{E}_{P_{\text{data}}(\mathbf{x}, \mathbf{y})} [d(\mathbf{x}, \mathbf{y}; \phi)] - \mathbb{E}_{\mathbf{z} \sim P_Z; \mathbf{y} \sim P_Y} [d(\mathbf{g}(\mathbf{z}, \mathbf{y}; \theta); \phi)].$$

Optimality

Optimal values (min-max problem):

$$(\theta^*, \phi^*) = \operatorname{argmin}_{\theta} \operatorname{argmax}_{\phi, d \in \text{Lip}} L(\theta, \phi).$$

With infinite data, infinite critic capacity achieving the maximum implies,

$$\theta^* = \operatorname{argmin}_{\theta} W_1(P_{\text{data}}, P_g P_Y).$$

If this minima is achieved $W_1(P_{\text{data}}, P_{g^*} P_Y) = 0$. This implies weak convergence of the measures. Which in turn means,

$$\mathbb{E}_{P_{\text{data}}} [I(x, y)] = \mathbb{E}_{z \sim P_Z, y \sim P_Y} [I(g(z, y; \theta^*), y)], \forall I \in C_b(\Omega_X \times \Omega_Y).$$

Optimality

Let $\hat{\mathbf{y}}$ be given for which $P_Y(\hat{\mathbf{y}}) \neq 0$. Let the true and approximated conditional expectations be bounded. Then given $\epsilon > 0$ and $q \in C_b(\Omega_X)$, there exists a positive integer $\tilde{N} := \tilde{N}(\hat{\mathbf{y}}, q, \epsilon)$ such that

$$\left| \mathbb{E}_{P_{\text{data}}(\mathbf{x}|\hat{\mathbf{y}})} [q(\mathbf{x})] - \mathbb{E}_{P_{g^*}(\mathbf{x}|\hat{\mathbf{y}})} [q(\mathbf{x})] \right| < \epsilon \quad \forall N_\theta \geq \tilde{N}$$

This implies

$$\mathbb{E}_{P_{\text{data}}(\mathbf{x}|\hat{\mathbf{y}})} [q(\mathbf{x})] \approx \mathbb{E}_{P_Z} [q(\mathbf{g}(\mathbf{z}, \hat{\mathbf{y}}; \boldsymbol{\theta}^*))].$$

Alternate Sampling Strategy

Under the same conditions as the previous theorem, given $\epsilon > 0$ and $q \in C_b(\Omega_X)$, there exists a positive real number $\bar{\sigma} := \bar{\sigma}(\hat{y}, q, \epsilon)$ and a positive integer with $\bar{N} := \bar{N}(\bar{\sigma})$ such that

$$\left| \mathbb{E}_{P_{\text{data}}(x|\hat{y})} [q(x)] - \mathbb{E}_{y \sim P_{Y_\sigma}, x \sim P_{g^*}(x|\hat{y})} [q(x)] \right| < \epsilon \quad \forall N_\theta \geq \bar{N}$$

Here $P_{Y_\sigma} = \mathcal{N}(\hat{y}, \sigma I)$.

This implies

$$\mathbb{E}_{P_{\text{data}}(x|\hat{y})} [q(x)] \approx \mathbb{E}_{z \sim P_Z, y \sim P_{Y_\sigma}} [q(\mathbf{g}(z, y; \theta^*))].$$

Implementation of CWGAN

- In practice, the “loss” term is approximated as:

$$L(\theta, \phi) = \frac{1}{N} \left(\sum_{i=1}^N d(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}; \phi) - \sum_{i=1}^N d(\mathbf{g}(\mathbf{z}^{(i)}, \mathbf{y}^{(i)}; \theta), \mathbf{y}^{(i)}; \phi) \right) + \lambda P(\phi).$$

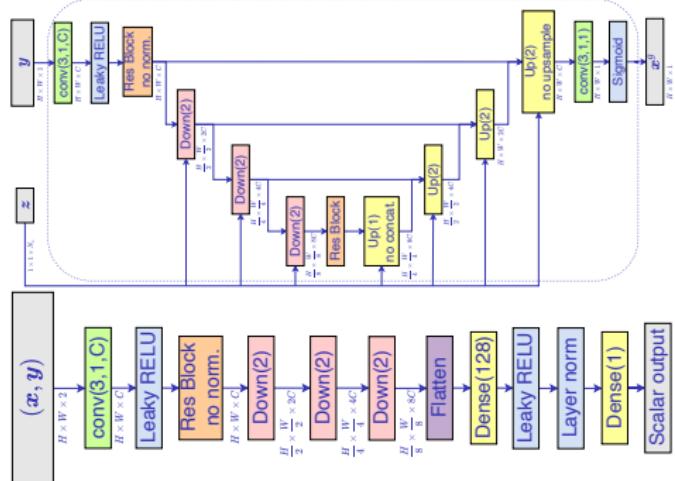
$$P(\phi) = \frac{1}{N} \sum_{i=1}^N (\|\nabla d(\mathbf{h}^{(i)}, \mathbf{y}^{(i)}; \phi)\|_2 - 1)^2$$

$$\mathbf{h}^{(i)} = \delta^{(i)} \mathbf{x}^{(i)} + (1 - \delta^{(i)}) \mathbf{g}(\mathbf{z}^{(i)}, \mathbf{y}^{(i)}; \theta), \quad \delta^{(i)} \sim U(0, 1).$$

- In the second term the gradient is wrt. both arguments.
- The min-max problem is hard to solve. Use stochastic gradient descent-ascent. Five steps of ascent followed by one descent.
- Hard to track convergence. Look for stability of the “loss” function and track statistics of generated samples. For a given \mathbf{y} compute the generated \mathbf{x} ’s and compare with the true \mathbf{x} .

Implementation of CWGAN

- If \mathbf{X} and \mathbf{Y} are vectors, the generator and the critic take the form of an MLP.
- If \mathbf{X} and \mathbf{Y} are images (discrete fields), the generator is a U-Net. Further \mathbf{z} is used to perform conditional instance normalization. The critic is a CNN - like a classifier.



Example 1

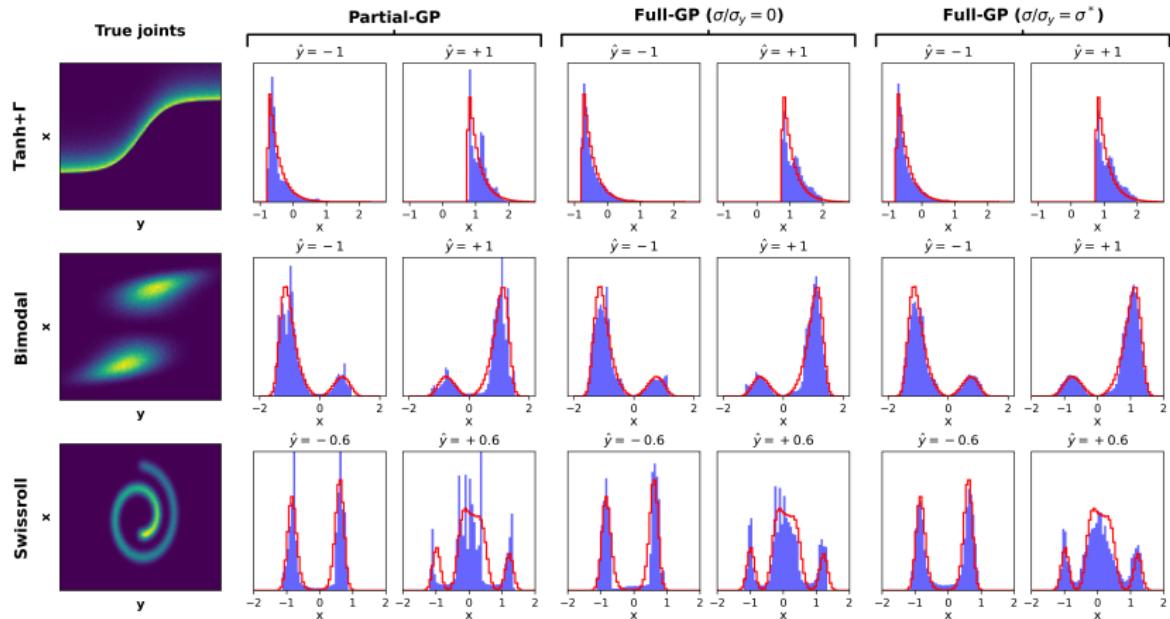


Figure: Histograms of 10^5 samples from $P_{g^*}(x|\hat{y})$. Red: outline the true conditional dist.

Example 2: Inverse heat equation

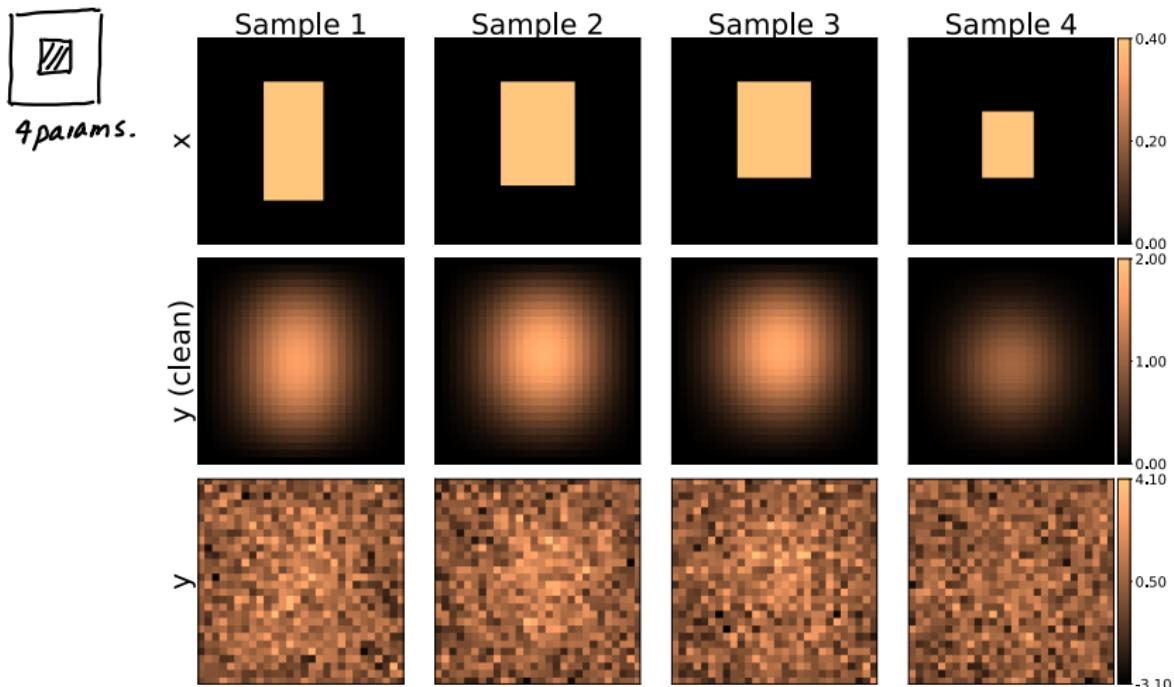


Figure: Samples from the rectangular prior dataset.

Example 2: Inverse heat equation

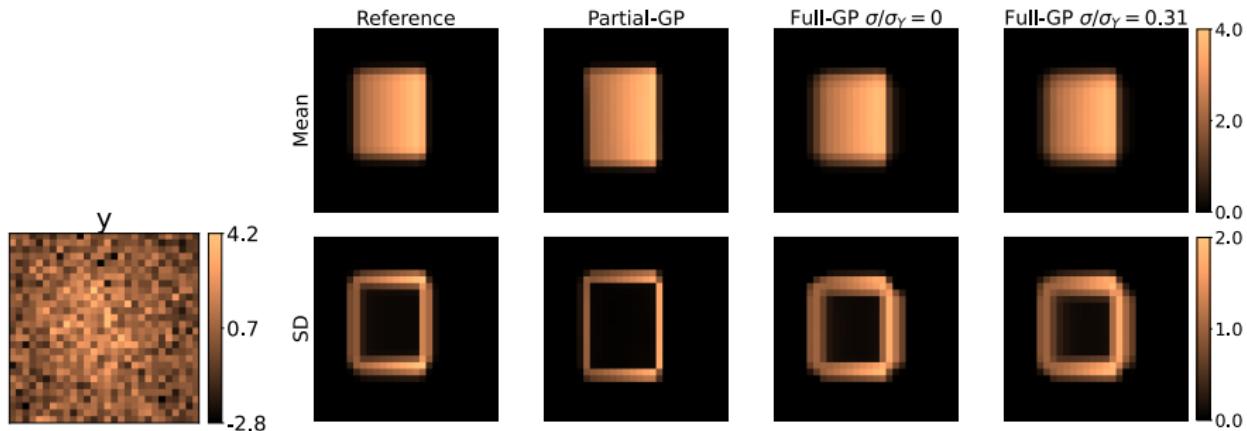


Figure: Measurement, Mean and SD from models and reference

Example 2: Inverse heat equation

Table: L^2 error in posterior statistics

Posterior Statistics	Partial-GP	Full-GP	
		$\sigma/\sigma_Y = 0$	$\sigma/\sigma_Y = 0.31$
Mean	0.441	0.420	0.406
SD	0.460	0.454	0.435

Example 3: Inverse Helmholtz equation

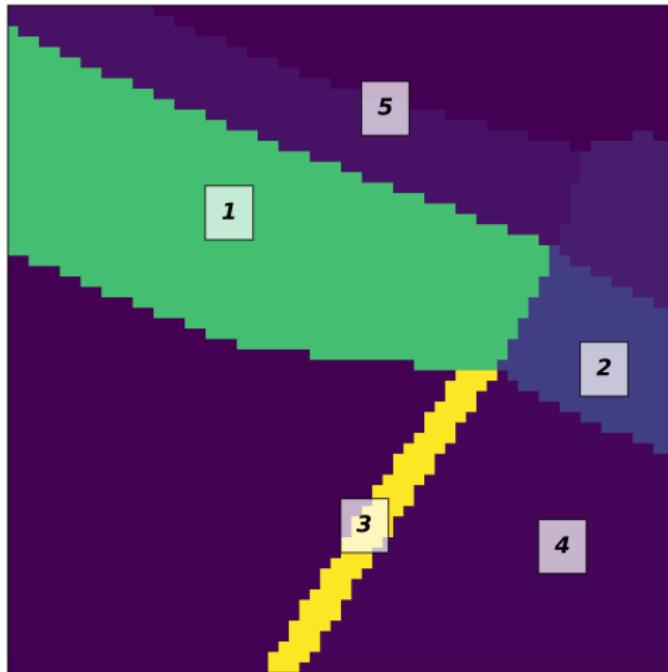


Figure: Geometry of the optical nerve head.

Example 3: Inverse Helmholtz equation

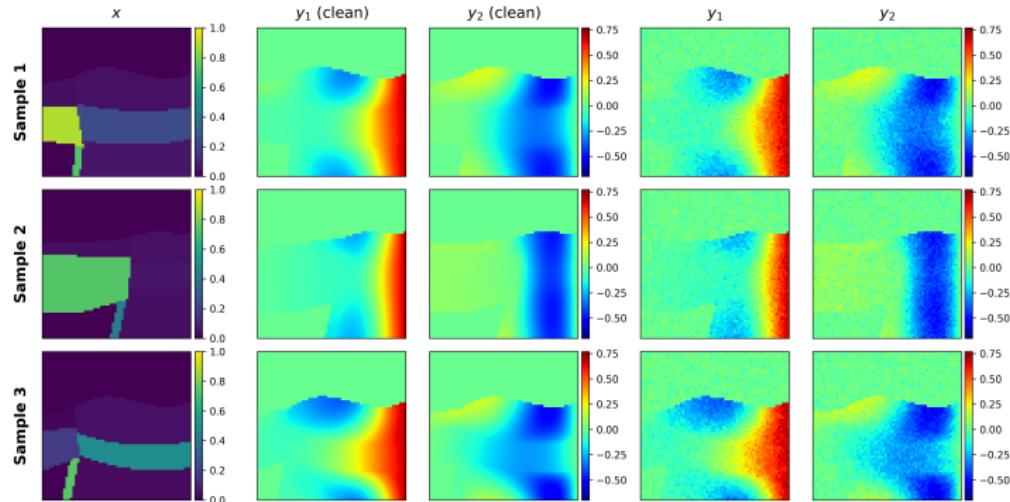


Figure: Training samples for the inverse Helmholtz problem.

Example 3: Inverse Helmholtz equation

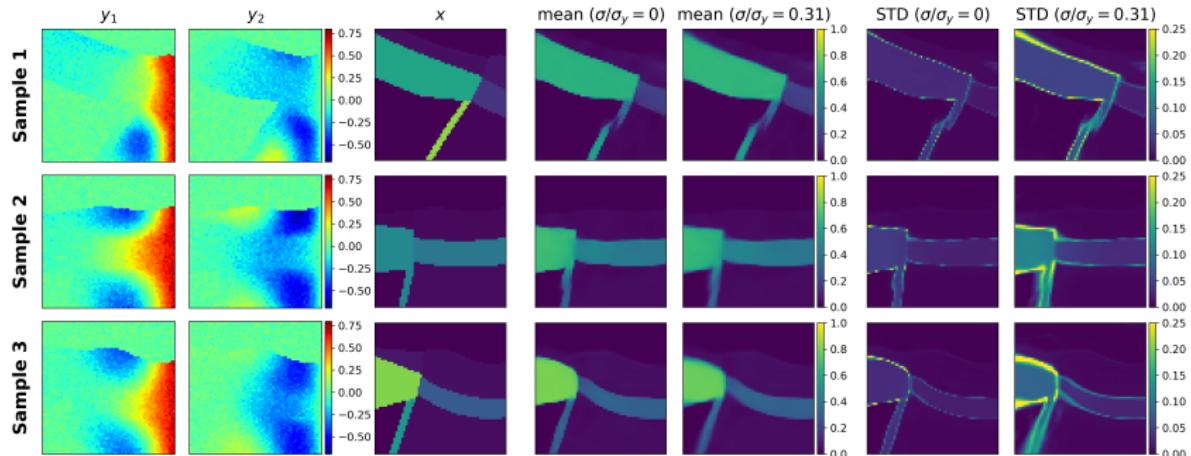


Figure: Point-wise mean and standard deviation.

Example 4: Inverse elasticity for tumor spheroids

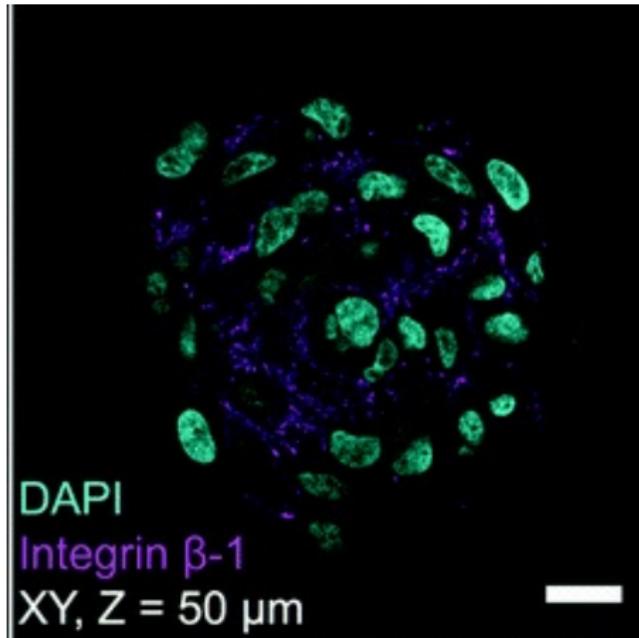


Figure: Confocal microscopy image of a tumor spheroid.

Example 4: Inverse elasticity for tumor spheroids (Kennedy & Foo)

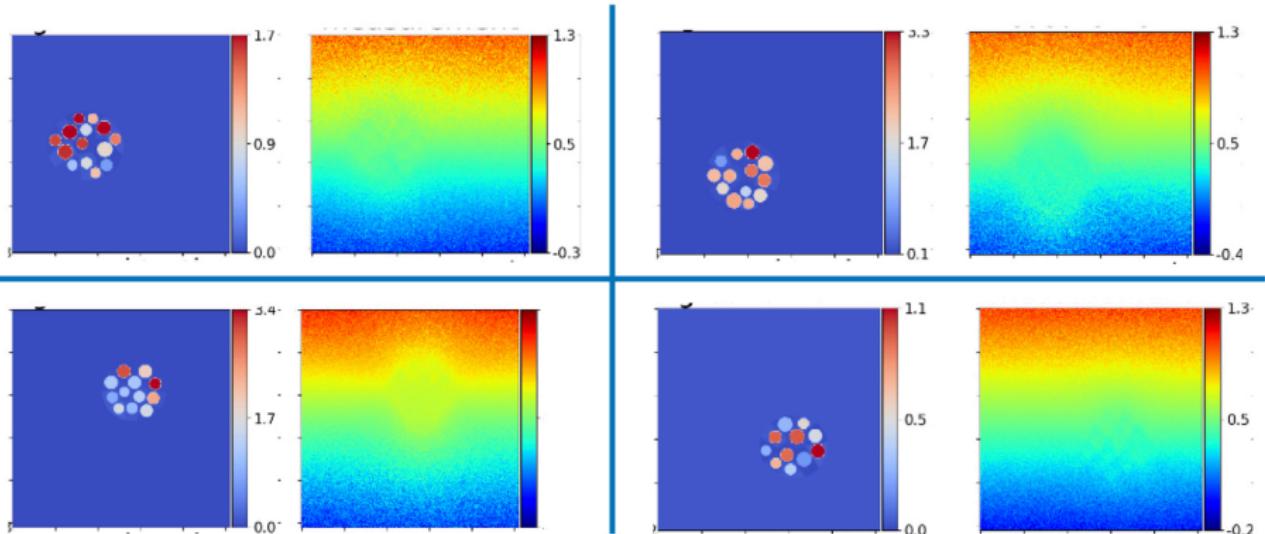


Figure: Training samples.

Example 4: Inverse elasticity for tumor spheroids

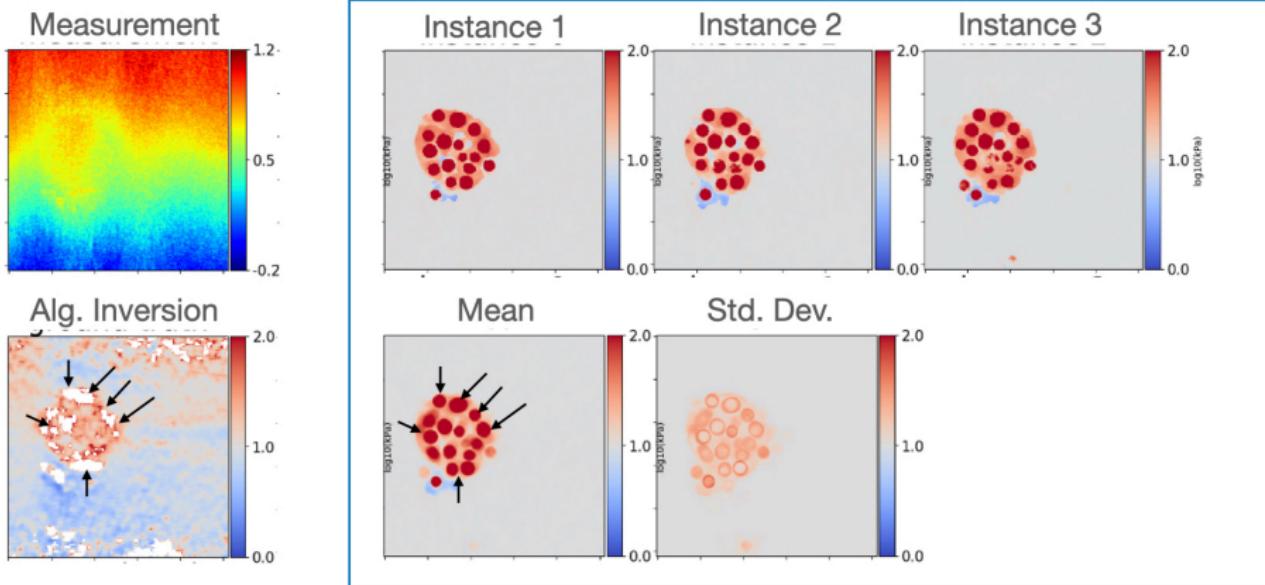
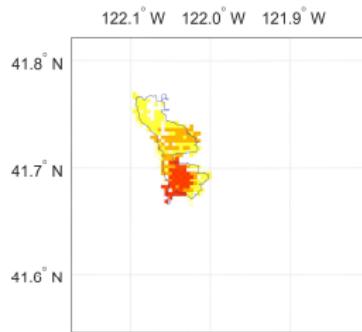
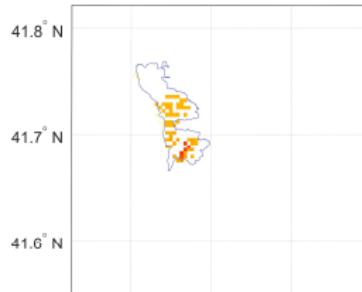


Figure: Generated samples, point-wise mean and standard deviation.

Example 5: Early Trajectory of a Wildfire (Shaddy, Mandel, Farguel..)



(b) Tenant

Figure: Satellite measurements of active fire regions

Example 5: Early Trajectory of a Wildfire

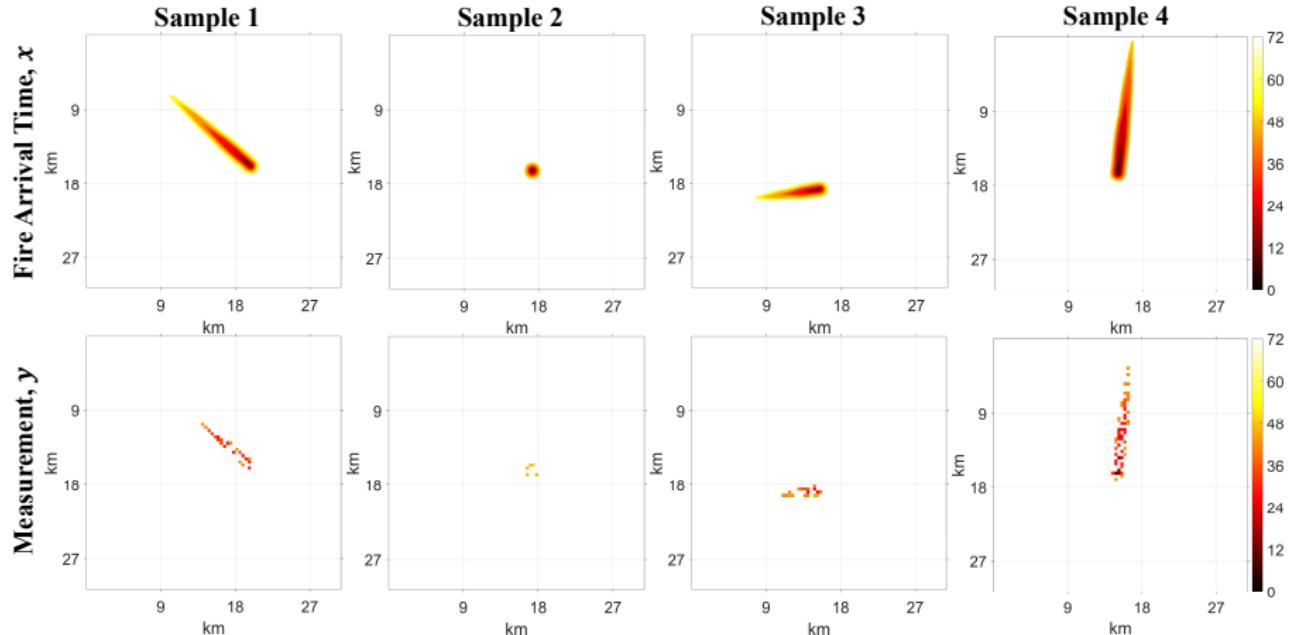
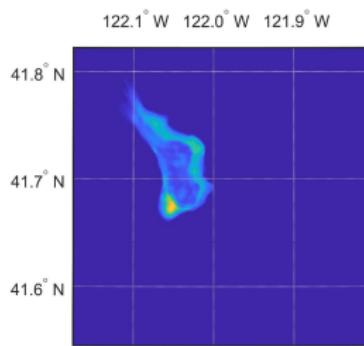
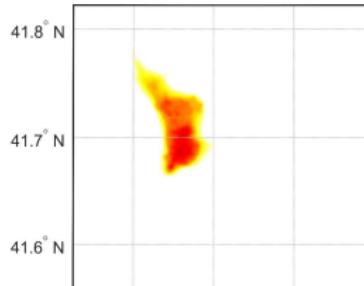


Figure: Training samples generated from WRF-SFIRE.

Example 5: Early Trajectory of a Wildfire



(b) Tenant

Figure: Fire arrival time map - mean and standard deviation.

Example 5: Early Trajectory of a Wildfire

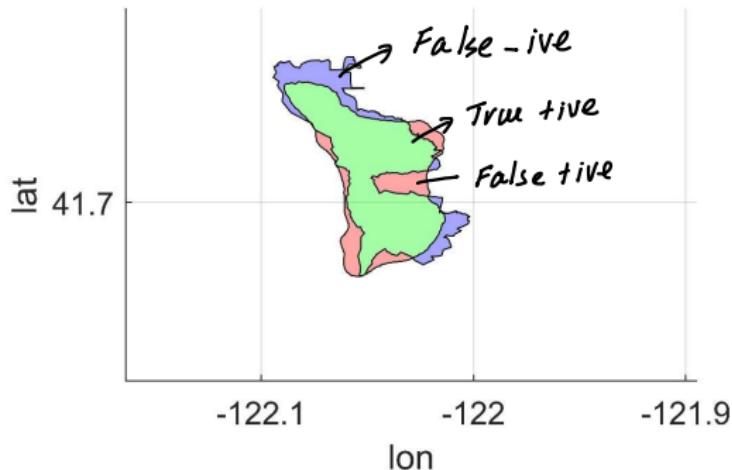


Figure: Validation against IR perimeter.

Table of Contents

Introduction

Wasserstein Generative Adversarial Network

Conditional Wasserstein Generative Adversarial Network

Generative Diffusion Networks

Conditional Generative Diffusion Networks

Generative Diffusion networks

Key references:

- Song, Yang, and Stefano Ermon. "Generative modeling by estimating gradients of the data distribution." *Advances in neural information processing systems* 32 (2019).
- Hyvarinen, Aapo, and Peter Dayan. "Estimation of non-normalized statistical models by score matching." *Journal of Machine Learning Research* 6.4 (2005).
- Vincent, Pascal. "A connection between score matching and denoising autoencoders." *Neural computation* 23.7 (2011): 1661-1674.
- Batzolis, Georgios, et al. "Conditional image generation with score-based diffusion models." *arXiv preprint arXiv:2111.13606* (2021).

Main idea

- Any attempt to approximate P_{data} has to ensure that the density must integrate to 1. Since N_x is large, this involves a challenging integral.
- On the other hand consider the score: $\mathbf{s}_{\text{data}} : \mathbb{R}^{N_x} \rightarrow \mathbb{R}^{N_x}$

$$\mathbf{s}_{\text{data}} = \nabla \log P_{\text{data}} = \nabla \log \left(\frac{f}{Z} \right) = \nabla (\log f - \log Z) = \frac{\nabla f}{f}.$$

does not depend on the normalizing constant.

- And Langevin equation offers a way to use the score to generate samples.
- Step 1: learn the score. Step 2: use Langevin equation to generate samples.

Step 1: Learning the score

- Let $s(x; \theta)$ be an approximation to $s_{\text{data}}(x)$.
- To find θ minimize

$$L(\theta) = \frac{1}{2} \int_{\Omega_X} |s(x; \theta) - s_{\text{data}}(x)|^2 P_{\text{data}}(x) dx.$$

- Cannot work with this directly since s_{data} is not known. Work on this in the next page.....

Step 1: Learning the score

$$\begin{aligned} L(\theta) &= \frac{1}{2} \int_{\Omega_X} |\mathbf{s}(\mathbf{x}; \theta) - \mathbf{s}_{\text{data}}(\mathbf{x})|^2 P_{\text{data}}(\mathbf{x}) d\mathbf{x} \\ &= \frac{1}{2} \int_{\Omega_X} |\mathbf{s}|^2 P_{\text{data}}(\mathbf{x}) d\mathbf{x} - \int_{\Omega_X} \mathbf{s} \cdot \mathbf{s}_{\text{data}} P_{\text{data}}(\mathbf{x}) d\mathbf{x} \\ &\quad + \frac{1}{2} \int_{\Omega_X} |\mathbf{s}_{\text{data}}|^2 P_{\text{data}}(\mathbf{x}) d\mathbf{x}. \\ &= \frac{1}{2} \int_{\Omega_X} |\mathbf{s}|^2 P_{\text{data}}(\mathbf{x}) d\mathbf{x} - \int_{\Omega_X} \mathbf{s} \cdot \frac{\nabla P_{\text{data}}}{P_{\text{data}}} P_{\text{data}}(\mathbf{x}) d\mathbf{x}. \\ &= \frac{1}{2} \int_{\Omega_X} |\mathbf{s}|^2 P_{\text{data}} d\mathbf{x} + \int_{\Omega_X} \nabla \cdot \mathbf{s} P_{\text{data}} d\mathbf{x} - \int_{\partial \Omega_X} \mathbf{s} \cdot \mathbf{n} P_{\text{data}} d\mathbf{x}. \\ &= \mathbb{E}_{P_{\text{data}}} \left[\frac{1}{2} |\mathbf{s}|^2 + \nabla \cdot \mathbf{s} \right]. \end{aligned}$$

Step 1: Learning the score

So,

$$\begin{aligned} L(\theta) &= \mathbb{E}_{P_{\text{data}}} \left[\frac{1}{2} |\mathbf{s}|^2 + \nabla \cdot \mathbf{s} \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \left(\frac{1}{2} |\mathbf{s}(\mathbf{x}^{(i)}; \theta)|^2 + \sum_{j=1}^{N_x} \frac{\partial}{\partial x_j} s_j(\mathbf{x}^{(i)}; \theta) \right). \end{aligned}$$

- This will approximate the score, however, the approximation will be poor in regions where P_{data} is small.
- We will see in the next step we need a good approximation everywhere.

Step 2: Using the score to generate

- Given $\mathbf{x}_0 \sim P_0(\mathbf{x})$, consider the iterations

$$\mathbf{x}_t = \mathbf{x}_{t-1} + \frac{\epsilon}{2} \nabla \log P_{\text{data}}(\mathbf{x}_{t-1}) + \sqrt{\epsilon} \mathbf{z}_t,$$

where $\mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$, and $\epsilon \rightarrow 0$ as $t \rightarrow \infty$, then in the limit $\mathbf{x}_t \sim P_{\text{data}}$.

- This equation is the Euler-Maruyama discretization of the Ito ODE

$$d\mathbf{X}_t = \frac{\epsilon}{2} \nabla \log P_{\text{data}}(\mathbf{x}) + \sqrt{\epsilon} \mathbf{z}_t,$$

with $\mathbf{x}_0 \sim P_0(\mathbf{x})$.

Step 2: Using the score to generate

- The corresponding Fokker-Plank equation governs the evolution of the probability density of \mathbf{X}_t and is given by,

$$\begin{aligned}\frac{\partial P_X(\mathbf{x}, t)}{\partial t} = & -\nabla \cdot \left(\frac{\epsilon}{2} \nabla \log(P_{\text{data}}(\mathbf{x})) P_X(\mathbf{x}, t) \right) \\ & + \nabla \cdot \nabla \cdot \left(\frac{\epsilon}{2} \mathbf{1} P_X(\mathbf{x}, t) \right),\end{aligned}$$

with $P_X(\mathbf{x}, 0) \sim P_0(\mathbf{x})$.

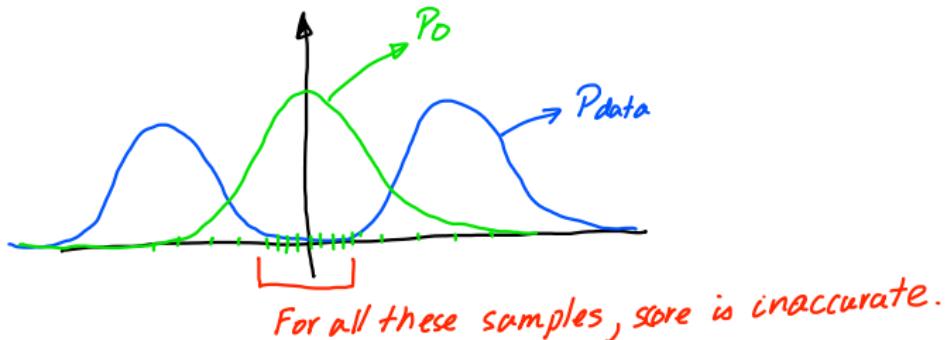
- The steady state solution is given by

$$\nabla \cdot \left(-\frac{\nabla P_{\text{data}}(\mathbf{x})}{P_{\text{data}}(\mathbf{x})} P_X(\mathbf{x}, \infty) + \nabla P_X(\mathbf{x}, \infty) \right) = 0.$$

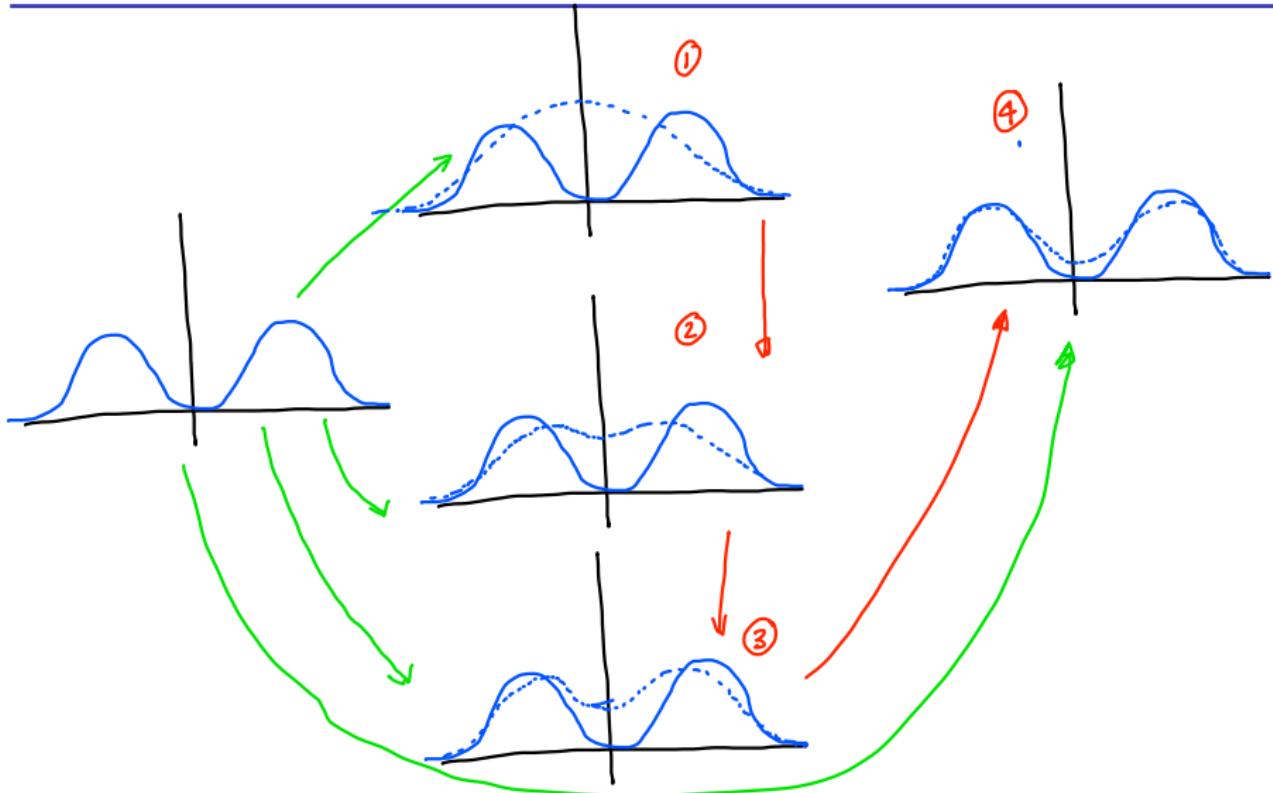
- A solution is $P_X(\mathbf{x}, \infty) = P_{\text{data}}(\mathbf{x})$. It is unique!
-

Step 2: Using the score to generate

This would work however the approximation of the real score will be poor in regions where P_{data} is small.



Remedy



Remedy

Initially work with a smoother version P_{data} that has support near P_0 .

- Consider

$$\begin{aligned} P_\sigma(\tilde{\mathbf{x}}) &= \int_{\Omega_X} P_{\text{data}}(\mathbf{x}) P_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) d\mathbf{x} \\ P_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) &= \exp\left(-\frac{|\tilde{\mathbf{x}} - \mathbf{x}|^2}{2\sigma^2}\right)/\mathbb{Z}. \end{aligned}$$

- Now we need a NN $\mathbf{s}(\tilde{\mathbf{x}}, \sigma; \theta)$ that will approximate

$$\mathbf{s}_\sigma(\tilde{\mathbf{x}}) = \nabla \log P_\sigma(\tilde{\mathbf{x}}).$$

- So we want to minimize

$$L(\theta, \sigma) = \frac{1}{2} \int_{\Omega_X} |\mathbf{s}(\tilde{\mathbf{x}}, \sigma; \theta) - \mathbf{s}_\sigma(\tilde{\mathbf{x}})|^2 P_\sigma(\tilde{\mathbf{x}}) d\tilde{\mathbf{x}}.$$

But have access only to samples from P_{data} .

Step 1: Learning the score

$$L(\theta, \sigma) = \int_{\Omega_X} \left(\frac{1}{2} |\mathbf{s}|^2 - \mathbf{s} \cdot \mathbf{s}_\sigma + \frac{1}{2} |\mathbf{s}_\sigma|^2 \right) P_\sigma(\tilde{\mathbf{x}}) d\tilde{\mathbf{x}}.$$

Consider the second term

$$\begin{aligned} \text{2nd term} &= \int_{\Omega_X} \mathbf{s} \cdot \mathbf{s}_\sigma P_\sigma(\tilde{\mathbf{x}}) d\tilde{\mathbf{x}} \\ &= \int_{\Omega_X} \mathbf{s} \cdot \nabla P_\sigma(\tilde{\mathbf{x}}) d\tilde{\mathbf{x}} \\ &= \int_{\Omega_X} \mathbf{s} \cdot \nabla \left(\int_{\Omega_X} P_{\text{data}}(\mathbf{x}) P_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) d\mathbf{x} \right) d\tilde{\mathbf{x}} \\ &= \int_{\Omega_X} \int_{\Omega_X} \mathbf{s} \cdot \nabla \left(\log P_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) \right) P_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) P_{\text{data}}(\mathbf{x}) d\mathbf{x} d\tilde{\mathbf{x}}. \end{aligned}$$

Step 1: Learning the score

Consider the first term

$$\begin{aligned}\text{1st term} &= \int_{\Omega_X} \frac{1}{2} |\mathbf{s}|^2 P_\sigma(\tilde{\mathbf{x}}) d\tilde{\mathbf{x}} \\ &= \int_{\Omega_X} \frac{1}{2} |\mathbf{s}|^2 \left(\int_{\Omega_X} P_{\text{data}}(\mathbf{x}) P_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) d\mathbf{x} \right) d\tilde{\mathbf{x}} \\ &= \int_{\Omega_X} \int_{\Omega_X} \frac{1}{2} |\mathbf{s}|^2 P_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) P_{\text{data}}(\mathbf{x}) d\mathbf{x} d\tilde{\mathbf{x}}.\end{aligned}$$

Step 1: Learning the score

Combine these two

$$\begin{aligned} L(\theta, \sigma) &= \int_{\Omega_X} \int_{\Omega_X} \left(\frac{1}{2} |\mathbf{s}|^2 - \mathbf{s} \cdot \nabla \log P_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) \right) P_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) P_{\text{data}}(\mathbf{x}) d\mathbf{x} d\tilde{\mathbf{x}} \\ &= \frac{1}{2} \int_{\Omega_X} \int_{\Omega_X} |\mathbf{s} - \nabla \log P_\sigma(\tilde{\mathbf{x}}|\mathbf{x})|^2 P_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) P_{\text{data}}(\mathbf{x}) d\mathbf{x} d\tilde{\mathbf{x}} \\ &= \frac{1}{2} \int_{\Omega_X} \int_{\Omega_X} |\mathbf{s} + \frac{1}{\sigma^2} (\tilde{\mathbf{x}} - \mathbf{x})|^2 P_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) P_{\text{data}}(\mathbf{x}) d\mathbf{x} d\tilde{\mathbf{x}} \\ &\approx \frac{1}{2N} \sum_{i=1}^N |\mathbf{s}(\tilde{\mathbf{x}}^{(i)}, \sigma; \theta) + \frac{1}{\sigma^2} (\tilde{\mathbf{x}}^{(i)} - \mathbf{x}^{(i)})|^2 \end{aligned}$$

Step 1: Learning the score

We will train $s(\tilde{\mathbf{x}}, \sigma; \theta)$ for large and small values of σ . Therefore select $\sigma_j \in (\sigma_1, \dots, \sigma_L)$ with $\frac{\sigma_{j+1}}{\sigma_j} = \rho > 1$. And then overall loss function for training the score-matching NN is

$$L(\theta) = \frac{1}{2LN} \sum_{j=1}^L \sum_{i=1}^N |s(\tilde{\mathbf{x}}^{(i)}, \sigma_j; \theta) + \frac{1}{\sigma_j^2} (\tilde{\mathbf{x}}^{(i)} - \mathbf{x}^{(i)})|^2$$

Step 2: Using the score to generate

Annealed Langevin equation:

- $\tilde{\mathbf{x}} \sim P_0$
- For $j = 1, \dots, L$
- $\alpha_j = \epsilon \rho^{L-j}$
- For $t = 1, \dots, T$
 - $\tilde{\mathbf{x}}_t = \tilde{\mathbf{x}}_{t-1} + \frac{\alpha_j}{2} \mathbf{s}(\tilde{\mathbf{x}}_{t-1}, \sigma_j; \theta^*) + \sqrt{\alpha_j} \mathbf{z}_t$
- End (t)
- End (j)

Table of Contents

Introduction

Wasserstein Generative Adversarial Network

Conditional Wasserstein Generative Adversarial Network

Generative Diffusion Networks

Conditional Generative Diffusion Networks

Main idea: in generative diffusion networks if we replace $P_{\text{data}}(\mathbf{x})$ by $P_{\text{data}}(\mathbf{x}|\mathbf{y})$, then things would be fine. Let us re-examine Step 1 and Step 2 keeping this in mind.

Step 1: Learning the score

- Consider

$$P_\sigma(\tilde{\mathbf{x}}|\hat{\mathbf{y}}) = \int_{\Omega_X} P_{\text{data}}(\mathbf{x}|\mathbf{y}) P_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) d\mathbf{x}$$

- Now we need a NN $\mathbf{s}(\tilde{\mathbf{x}}, \sigma, \mathbf{y}; \theta)$ that will approximate

$$\mathbf{s}_\sigma(\tilde{\mathbf{x}}|\mathbf{y}) = \nabla \log P_\sigma(\tilde{\mathbf{x}}|\mathbf{y}).$$

- So we want to minimize

$$L(\theta, \sigma) = \frac{1}{2} \int_{\Omega_Y} \int_{\Omega_X} |\mathbf{s}(\tilde{\mathbf{x}}, \sigma, \mathbf{y}; \theta) - \mathbf{s}_\sigma(\tilde{\mathbf{x}}|\mathbf{y})|^2 P_\sigma(\tilde{\mathbf{x}}) d\tilde{\mathbf{x}} P_Y(\mathbf{y}) d\mathbf{y}.$$

But have access only to samples from $P_{\text{data}}(\mathbf{x}, \mathbf{y})$.

Step 1: Learning the score

- Following the same steps....

$$\begin{aligned} L(\theta, \sigma) &= \frac{1}{2} \int_{\Omega_Y} \int_{\Omega_X} \int_{\Omega_X} |\mathbf{s} + \frac{1}{\sigma^2}(\tilde{\mathbf{x}} - \mathbf{x})|^2 \\ &\quad P_\sigma(\tilde{\mathbf{x}}|\mathbf{x})P_{\text{data}}(\mathbf{x}|\mathbf{y})P_Y(\mathbf{y})d\mathbf{x}d\tilde{\mathbf{x}}d\mathbf{y} \\ &= \int_{\Omega_X} \int_{\Omega_Y} \int_{\Omega_X} |\mathbf{s} + \frac{1}{\sigma^2}(\tilde{\mathbf{x}} - \mathbf{x})|^2 \\ &\quad P_\sigma(\tilde{\mathbf{x}}|\mathbf{x})P_{\text{data}}(\mathbf{x}, \mathbf{y})d\mathbf{x}d\mathbf{y}d\tilde{\mathbf{x}} \\ &\approx \frac{1}{2N} \sum_{i=1}^N |\mathbf{s}(\tilde{\mathbf{x}}^{(i)}, \sigma, \mathbf{y}^{(i)}; \theta) + \frac{1}{\sigma^2}(\tilde{\mathbf{x}}^{(i)} - \mathbf{x}^{(i)})|^2. \end{aligned}$$

- Step 2 remains unchanged.