

# A Deeper Understanding of the Quadratic Wasserstein Metric in Inverse Data Matching

---

**Yunan Yang** (NYU), Matt Dunlop (NYU), Björn Engquist (UT-Austin), Kui Ren (Columbia University)

May 5, 2020

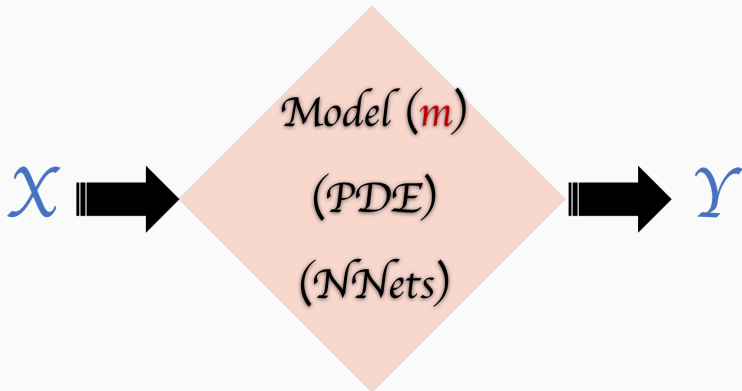
<https://arxiv.org/abs/1911.06911>

<https://arxiv.org/abs/2004.03730>

This work is partially supported by NSF DMS-1913129.

MSRI Workshop: Optimal Transport And Applications To Machine Learning And Statistics, May 4, 2020 - May 8, 2020

# Inverse Data Matching Problem



Inverse data matching problems aim at finding  $m$  such that the predicted outputs  $(X, Y(m))$  match given measured data  $(X, Y)$ .

## The Role of the Wasserstein Metric

- As an objective function measuring data,  $W_p(Y(m), Y)$ .
- The functional space for  $m$  (new gradient formula).
- Study the convergence of  $m_k$  as iteration number  $k$  increases (gradient flow, JKO, etc).
- Study the convergence of  $m_n$  (as an empirical measure) as (over)parameterization  $n$  increases (mean-field limit).
- etc.

# The Role of the Wasserstein Metric

- As an objective function measuring data,  $W_p(Y(m), Y)$ .
- The functional space for  $m$  (new gradient formula).
- Study the convergence of  $m_k$  as iteration number  $k$  increases (gradient flow, JKO, etc).
- Study the convergence of  $m_n$  (as an empirical measure) as parameterization  $n$  increases (mean-field limit).
- etc.

## The Model $F(m)$

Model ( $m$ )

(PDE)

(NNets)

$F$  is given; we just find  $m$  (e.g., PDEs).

OR

$F$  is not known;  $m$  depends on  $F$ .

## The Model $F(m)$

$F$  is given; we just find  $m$  (e.g., PDEs).

- Pro: We know the best (exact) forward problem!
- Con: The forward and inverse problems are so nonlinear!

OR

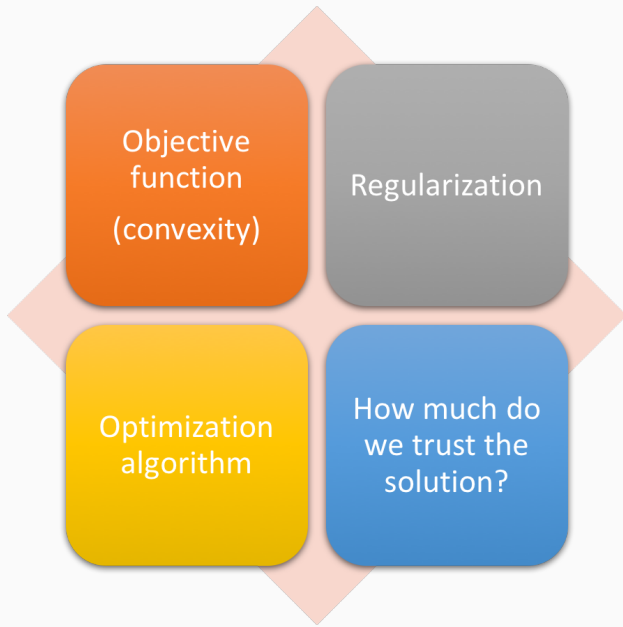
$F$  is not known; we are free to choose (e.g., XXX-net).

- Pro: The freedom to modify it to a “better” map (Over-Parametrization, ReLu)
- Con: Trial and error to build the model

# **1. Better Convexity (Optimization Landscape)**

---

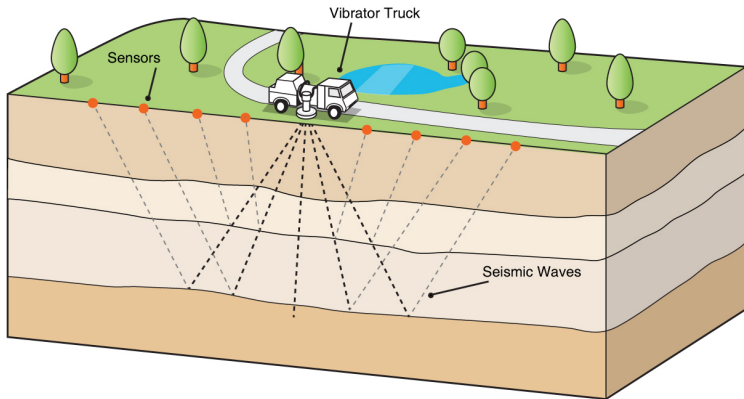
## Important Components in the Deterministic Approach



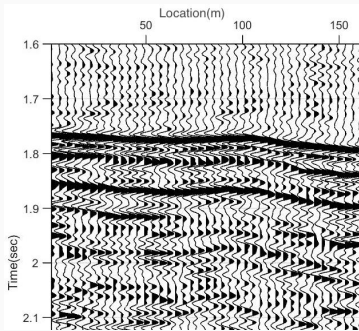


# Seismic Inversion: Earthquake Source, Hydrocarbons, etc.

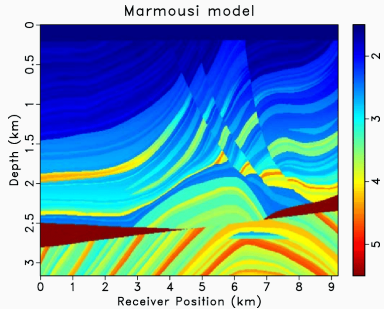
Seismic inversion is one of the inherently more difficult families of large-scale nonlinear inverse problems.



# Seismic Inversion



Invert  
⇒



Waveform measurements from receivers at the surface

Subsurface properties (i.e. wave velocity or material density)

## Forward Problem

$$\mathcal{F} : m \rightarrow u|_{\Gamma}, \Gamma \subseteq \partial\Omega \text{ or } \Omega$$

## Inverse Problem

$$\mathcal{G} : u|_{\Gamma} \rightarrow m$$

$\mathcal{F}$  and  $\mathcal{G}$  are often nonlinear.

$$m^* = \underset{m}{\operatorname{argmin}} J(f(m), g)$$

$$f(m) = u|_{\Gamma}$$

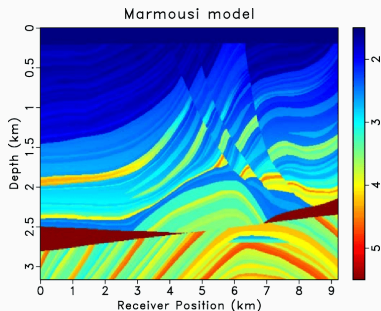
$J$  is an objective function measuring the difference between  $f$  and  $g$ .

# Seismic (Nonlinear) Inversion

## Forward Wave Propagation

$$\left\{ \begin{array}{l} m(\mathbf{x}) \frac{\partial^2 u(\mathbf{x}, t)}{\partial t^2} - \Delta u(\mathbf{x}, t) = s(\mathbf{x}, t) \\ \text{Zero i.c. in half-space } \Omega \\ \text{Neumann b.c. on } \partial\Omega \end{array} \right.$$

$$m(\mathbf{x}) = \frac{1}{c(\mathbf{x})^2}, \text{ } c(\mathbf{x}) \text{ is the wave velocity}$$



$m$

## Typical Effects from Variations in Wave Speed

The shift and dilation are typical effects from variations in velocity parameter  $m(x) = m$  (constant). For example:

$$\begin{cases} m \frac{\partial^2 u}{\partial t^2} = \frac{\partial^2 u}{\partial x^2}, & x > 0, t > 0, \\ u = 0, \quad \frac{\partial u}{\partial t} = 0, & x > 0, t = 0, \\ u = f(t), & x = 0, t > 0. \end{cases}$$

The solution to the equation is  $u(x, t; m) = f(t - \sqrt{mx})$ .

For fixed  $x$ , variation in  $m$  relates **shifts** in the signal.

For fixed  $t$ , variation in  $m$  generates the **dilation** in data.

## Traditional Least-Squares ( $L^2$ norm) Objective Function

$$J(m) = \frac{1}{2} \sum_r \int |f(x_r, t; m) - g(x_r, t)|^2 dt, \quad (1)$$

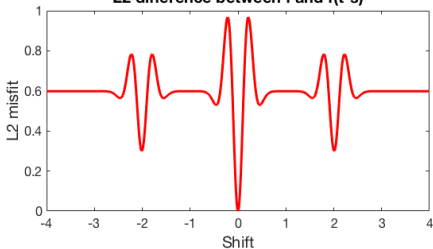
- observed data  $g$ ,
- simulated data  $f(m) = u|_{\Gamma}$ ,
- receiver  $x_r$ ,
- the model parameter  $m$ ,
- Regularization is often added in (1).

## Main Challenges

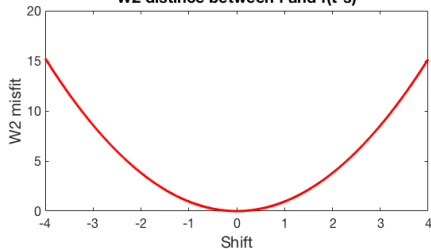
1. Local minima trapping
2. Sensitive to noise

# Motivation of Using the Wasserstein Distance (EMD)

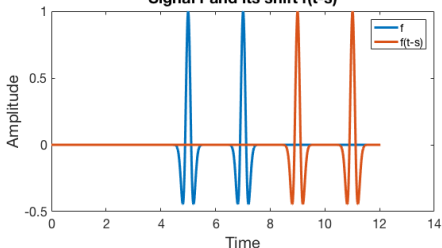
L2 difference between  $f$  and  $f(t-s)$



W2 distance between  $f$  and  $f(t-s)$



Signal  $f$  and its shift  $f(t-s)$



[Engquist-Froese, 2014] [Engquist, Froese & Y, 2016]

# The Quadratic Wasserstein Distance

## The Quadratic Wasserstein Distance

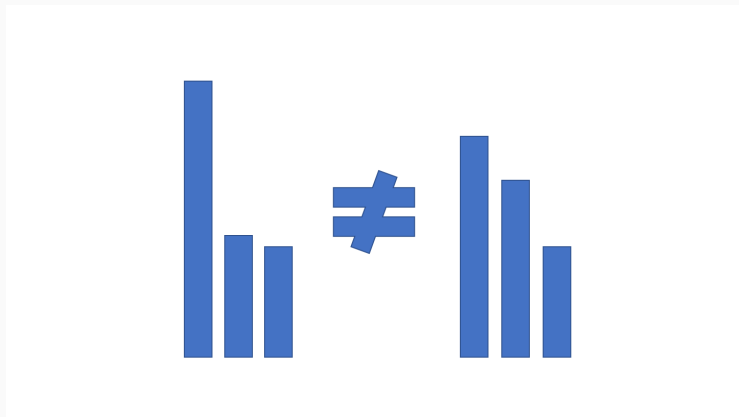
For  $f, g \in \mathcal{P}(\Omega)$  ( $f, g \geq 0$  and  $\int f = \int g = 1$ ), the quadratic Wasserstein distance is formulated as

$$W_2(f, g) = \left( \inf_{T \in \mathcal{M}} \int |x - T(x)|^p f(x) dx \right)^{\frac{1}{2}} \quad (2)$$

$\mathcal{M}$ : the set of all maps that rearrange the distribution  $f$  into  $g$ .



# Optimal Transport



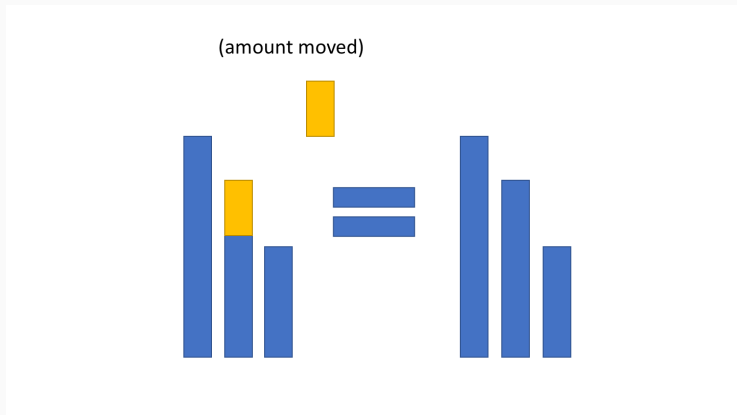
Synthetic data  $f$  (left) and observed data  $g$  (right)

# Optimal Transport



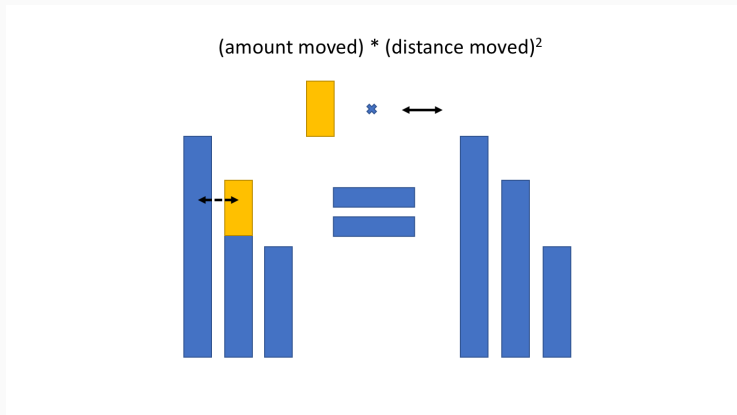
Synthetic data  $f$  (left) and observed data  $g$  (right)

# Optimal Transport



Synthetic data  $f$  (left) and observed data  $g$  (right)

# Optimal Transport



Synthetic data  $f$  (left) and observed data  $g$  (right)

# Tackling Nonconvexity

Let  $\{e_k\}_{k=1}^d$  be standard basis of the Euclidean space  $\mathbb{R}^d$ .

Assume  $s_k \in \mathbb{R}$ ,  $\lambda_k \in \mathbb{R}^+$ ,  $k = 1, \dots, d$  and  $A = \text{diag}(1/\lambda_1, \dots, 1/\lambda_d)$ .

We define  $f_\Theta$  as jointly the translation and dilation of  $g$ :

$$f_\Theta(x) = \det(A)g\left(A\left(x - \sum_{k=1}^d s_k e_k\right)\right), \Theta = \{s_1, \dots, s_d, \lambda_1, \dots, \lambda_d\}.$$

## Theorem (Convexity of $W_2$ in translation and dilation)

The optimal map between  $f_\Theta(x)$  and  $g(y)$  is  $y = T_\Theta(x)$  where

$$\langle T_\Theta(x), e_k \rangle = \frac{1}{\lambda_k}(\langle x, e_k \rangle - s_k), k = 1, \dots, d.$$

Moreover,  $I(\Theta) = W_2^2(f_\Theta(x), g)$  is a convex function of  $\Theta$ .

# Data Normalization: From Seismic Signal to Probability Density

- Absolute value scaling:  $f_2 = |f_1|$
- Square scaling:  $f_2 = f_1^2$
- Linear scaling:  $f_2 = f_1 + a$
- Exponential scaling:  $f_2 = \exp(af_1)$
- Soft-Plus:  $f_2 = \log(\exp(af_1) + 1)$

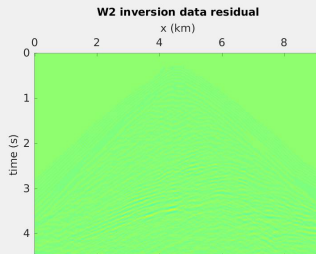
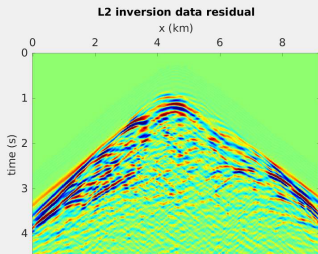
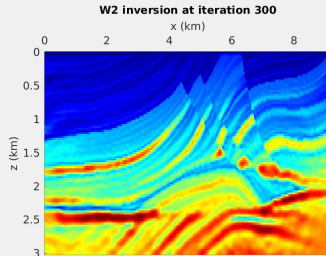
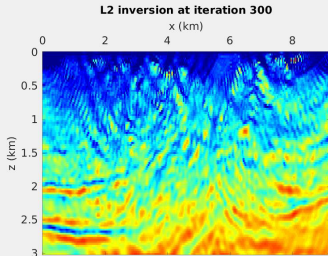
$$f = \frac{f_2}{\int f_2}$$

# Data Normalization: From Seismic Signal to Probability Density

- Absolute value scaling:  $f_2 = |f_1|$
- Square scaling:  $f_2 = f_1^2$
- Linear scaling:  $f_2 = f_1 + a$
- Exponential scaling:  $f_2 = \exp(af_1)$
- Soft-Plus:  $f_2 = \log(\exp(af_1) + 1)$

$$f = \frac{f_2}{\int f_2}$$

# Tackling Nonconvexity



[Y-Engquist-Sun-Hamfeldt, 2016]



## Convexity is not limited to the Wave Applications

- Inversion for Transport in homogeneous flow;
- Reconstruction from projections;
- Deconvolution of highly localized sources;
- Deconvolution from diffusive environment.

## **2. Robustness w.r.t. Noise**

---

## 2. More Robust w.r.t. Noise (Perturbation)

Given strictly positive probability density  $f = d\nu$ , we can define a Laplace-type linear operator

$$L = -\Delta + \nabla(-\log f) \cdot \nabla$$

which satisfies the fundamental integration by parts formula:

$$\int_{\mathbb{R}^d} (Lh_1)h_2 d\nu = \int_{\mathbb{R}^d} h_1(Lh_2) d\nu = \int_{\mathbb{R}^d} \nabla h_1 \cdot \nabla h_2 d\nu.$$

$$\|h\|_{L^2(f)}^2 = \int_{\mathbb{R}^d} h^2 d\nu, \quad \|h\|_{\dot{H}^1(f)}^2 = \int_{\mathbb{R}^d} |\nabla h|^2 d\nu,$$

$$\|h\|_{\dot{H}^{-1}(f)}^2 := \sup \left\{ \int_{\mathbb{R}^d} h\varphi d\nu \mid \|\varphi\|_{\dot{H}^1(f)}^2 \leq 1 \right\} = \int_{\mathbb{R}^d} h(L^{-1}h) d\nu.$$

If  $f = 1$ , we reconstruct the unweighted  $\mathcal{H}_{(\mathbb{R}^d)}^{-1}$  seminorm.

# The Connection with the Weak Norm

## Asymptotic Connection [Otto-Villani, 2000]

If  $\mu$  is the probability measure and  $d\pi$  is an infinitesimal perturbation that has zero total mass, then

$$W_2(\mu, \mu + d\pi) = \|d\pi\|_{\dot{\mathcal{H}}_{(d\mu)}^{-1}} + o(d\pi). \quad (3)$$

## Non-Asymptotic Connection [R. Peyre, 2018]

If both  $f = d\mu$  and  $g = d\nu$  are bounded from below and above by constants  $c_1$  and  $c_2$ , we have the following *non-asymptotic* equivalence between  $W_2$  and  $\dot{\mathcal{H}}_{(d\mu)}^{-1}$ :

$$\frac{1}{c_2} \|\mu - \nu\|_{\dot{\mathcal{H}}_{(\mathbb{R}^d)}^{-1}} \leq W_2(\mu, \nu) \leq \frac{1}{c_1} \|\mu - \nu\|_{\dot{\mathcal{H}}_{(\mathbb{R}^d)}^{-1}}, \quad (4)$$

## $H^s$ -based inverse matching

A linear inverse problem of finding  $m$  from noisy data  $g_\delta$

$$Am = g_\delta. \quad (5)$$

$A$  (a smoothing operator) is diagonal in the Fourier domain:

$$\widehat{A}(\xi) \sim \langle \xi \rangle^{-\alpha}. \quad (6)$$

We seek the solution by minimizing the objective functional

$$\mathcal{O}_{\mathcal{H}^s}(m) \equiv \frac{1}{2} \|f(m) - g\|_{\mathcal{H}^s}^2 := \frac{1}{2} \int_{\mathbb{R}^d} \langle \xi \rangle^{2s} |\widehat{f}(m)(\xi) - \widehat{g}(\xi)|^2 d\xi, \quad (7)$$

# What do we lose?

If we can obtain the solution by direct solve (best-case scenario)

## Theorem

Let  $R_c$  an approximation to  $A^{-1}$  defined through its symbol:

$$\widehat{R}_c(\boldsymbol{\xi}) \sim \begin{cases} \langle \boldsymbol{\xi} \rangle^\alpha, & |\boldsymbol{\xi}| < \xi_c \\ 0, & |\boldsymbol{\xi}| > \xi_c \end{cases}.$$

Let  $\delta = \|g_\delta - g\|_{\mathcal{H}^s}$ ,  $m_\delta^c := R_c g_\delta$  as the minimizer of  $\Phi(m)_{\mathcal{H}^s}$ .

$$\|m - m_\delta^c\|_{L^2} \lesssim \|m\|_{\mathcal{H}^{\frac{\alpha+\beta-s}{\beta}}} \delta^{\frac{\beta}{\alpha+\beta-s}}. \quad (8)$$

Reconstruction based on  $\mathcal{H}^s$  has an optimal spatial resolution

$$\varepsilon \sim \delta^{\frac{1}{\alpha+\beta-s}}. \quad (9)$$

# What do we gain?

If the noise contains mainly the higher frequency components

The solution at frequency  $\xi$  is therefore

$$\hat{m}(\xi) = \left( \hat{A}^*(\xi) (\langle \xi \rangle^{2s} \hat{A}) \right)^{-1} \hat{A}^*(\xi) \left( \langle \xi \rangle^{2s} \hat{g}_\delta(\xi) \right).$$

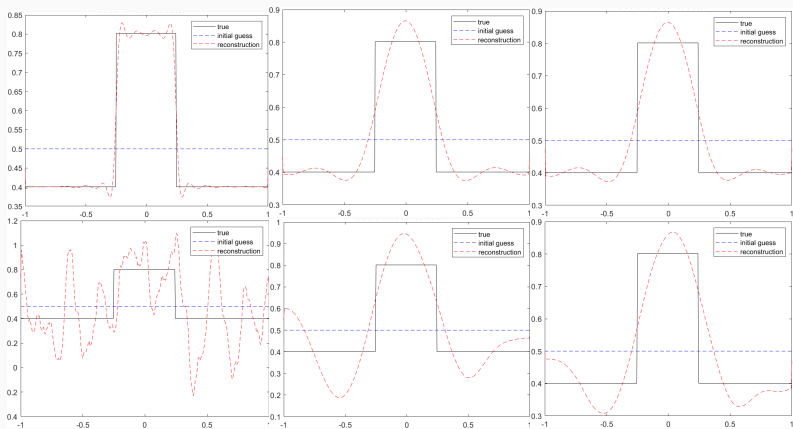
$$m = \left( A^* P A \right)^{-1} A^* P g_\delta, \quad P := (\mathcal{I} - \Delta)^{s/2},$$

where the operator  $(\mathcal{I} - \Delta)^{s/2}$  is defined through the relation

$$(\mathcal{I} - \Delta)^{s/2} m = \mathcal{F}^{-1} \left( \langle \xi \rangle^s \hat{m} \right),$$

$s = 0, s > 0, s < 0$ .

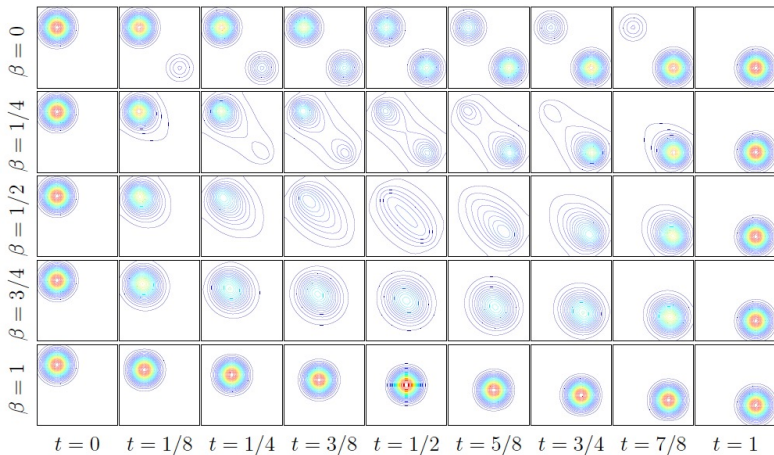
# What do we gain and loss?



Deconvolution with the kernel  $K_l(x) = \frac{1}{1+|x|}$  with the  $L^2$  (left),  $\mathcal{H}^{-1}$  (middle), and  $W_2$  (right) metrics. Top row: with noise-free data; Bottom row: with data containing respectively 2%, 10%, and 10% random noise.



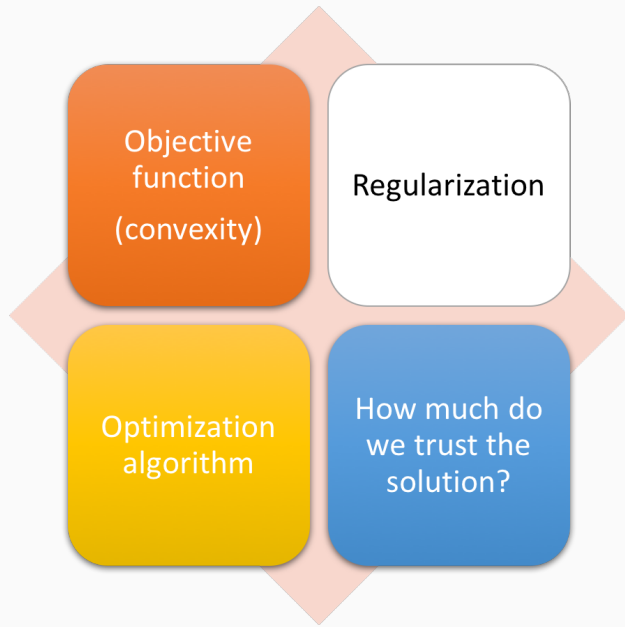
# Differences Between $W_2$ and $\dot{H}^{-1}$ (the gradient flow)



Top row: Geodesics in the  $\dot{H}^{-1}$  space

Bottom row: Geodesics in the  $W_2$  space

# Regularization



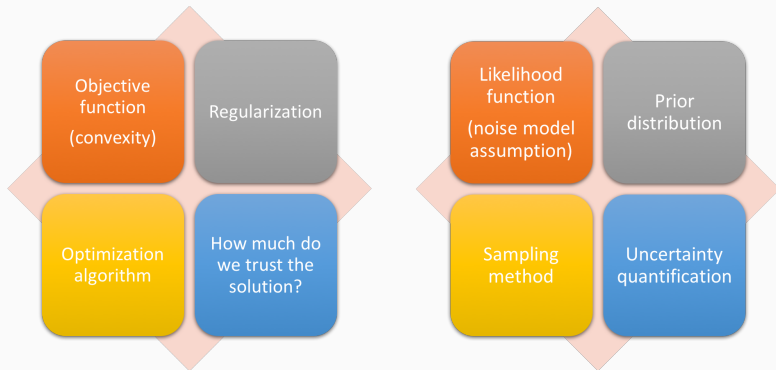
# Implicit Regularization

$$m^* = \underset{m}{\operatorname{argmin}} J(m) + R(m)$$

*Regularization does not have to be in the form of  $R(m)$ .*

- The choice of the objective function
- The choice of the data  
e.g., low-frequency data recovers low-wavenumber model
- The choice of numerical discretization
- The optimization algorithm (fixed step size)

# Deterministic & Bayesian



For **Large-Scale** inverse data matching problems

### **3. Wasserstein Metric as a Likelihood Function in Bayesian Inference**

---

# Data Normalization

**One problem:**  $\mathcal{G}(u)$  and  $y$  are not probability density functions.

**An potential solution:** Data Normalization; [Engquist-Y, 2020].

Given a  $\sigma : \mathbb{R} \rightarrow \mathbb{R}^+$ , we define  $P_\sigma$  on functions  $y : D \times T \rightarrow \mathbb{R}$  as

$$\tilde{y} = (P_\sigma y)(x, t) = \frac{1}{Z_\sigma(x)} \sigma(y(x, t)), \quad Z_\sigma(x) = \int_T \sigma(y(x, t')) dt'.$$

We only measure the  $T$  domain under the Wasserstein metric.

## $W_2$ Likelihood Function

$$W_2 \left( \widetilde{\mathcal{G}(u)}(x, \cdot), \widetilde{y}(x, \cdot) \right)^2 \approx \left\| \frac{\widetilde{\mathcal{G}(u)}(x, \cdot) - \widetilde{y}(x, \cdot)}{\widetilde{\mathcal{G}(u)}(x, \cdot)} \right\|_{\dot{H}^{-1}(\widetilde{\mathcal{G}(u)})}^2$$

which indicates the following noise model

$$\widetilde{y} = \eta \cdot \widetilde{\mathcal{G}(u)}, \quad \eta|u \sim N(1, \mathcal{L}(u))$$

where  $\mathcal{L}(u) : D(\mathcal{L}(u)) \rightarrow L^2(D; L^2(T))$  is defined by

$$\mathcal{L}(u)\varphi = - \underbrace{\frac{1}{\widetilde{\mathcal{G}(u)}}}_{\rho} \nabla_T \cdot \left( \underbrace{\widetilde{\mathcal{G}(u)}}_{\rho} \nabla_T \varphi \right),$$

where  $D(\mathcal{L}(u)) = \left\{ \varphi \in L^2(D; H^2(T)) \mid \int_T \varphi(\widetilde{\mathcal{G}(u)}) dt = 0 \right\}$  and  $\nabla_T$  is the gradient in the  $T$  domain.

# Likelihood Function

$\phi$	Likelihood function	Noise model assumption
$\Phi_{L^2}$	$\ \mathcal{G}(u)(x, \cdot) - y(x, \cdot)\ _{L^2(T)}^2$	$y = \mathcal{G}(u) + \eta, \eta \sim N(0, I)$
$\Phi_{H^{-1}}$	$\ \mathcal{G}(u)(x, \cdot) - y(x, \cdot)\ _{H^{-1}(T)}^2$	$y = \mathcal{G}(u) + \eta, \eta \sim N(0, -\Delta_T)$
$\Phi_{W_2}$	$W_2^2 \left( \widetilde{\mathcal{G}(u)}(x, \cdot), \widetilde{y}(x, \cdot) \right)$	$\widetilde{y} = \eta \cdot \widetilde{\mathcal{G}(u)}, \eta u \sim N(1, \mathcal{L}(u))$
$\Phi_M$	$\left\  \frac{\mathcal{G}(u)(x, \cdot) - y(x, \cdot)}{(y)(x, \cdot)} \right\ _{L^2(T)}^2$	$y = \eta \cdot \mathcal{G}(u), 1/\eta \sim N(1, I)$

The  $W_2$  metric can be regarded as asymptotically coming from **the state-dependent multiplicative noise data model:**  
*measurement error is proportional to the size of the quantity,  
and the distribution depends on the model parameter.*



### Theorem (Existence)

Let  $\pi_0$  be a Borel probability measure on  $X$ . Then for any choice  $\Phi \in \{\Phi_{L^2}, \Phi_{H^{-1}}, \Phi_{W_2}, \Phi_M\}$ ,

$$Z_\Phi(y) = \int_X \exp(-\Phi(u; y)) \pi_0(du)$$

is strictly positive and finite, and

$$\pi_\Phi^y(du) := \frac{1}{Z_\Phi(y)} \exp(-\Phi(u; y)) \pi_0(du)$$

defines a Radon probability measure on  $X$ .

### Theorem (Well-posedness)

Choose any  $\Phi \in \{\Phi_{L^2}, \Phi_{H^{-1}}, \Phi_{W_2}, \Phi_M\}$ . Under mild assumptions, there exists  $C_\Phi(r) > 0$  such that for all  $y, y' \in Y$  with

$$\|y\|_{L^\infty(D;L^\infty(T))}, \|y'\|_{L^\infty(D;L^\infty(T))} < r,$$

$$d_H(\pi_\Phi^y, \pi_\Phi^{y'}) \leq C_\Phi(r) \|y - y'\|_Y.$$

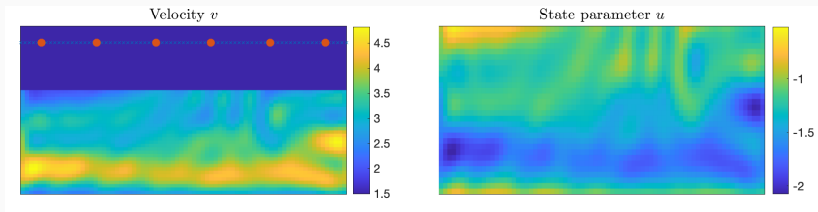
$d_H$  represents the Hellinger distance.

$$d_H(\pi_{\Phi_{W_2}}^y, \pi_{\Phi_{W_2}}^{y'}) \leq C_{W_2} \|y - y'\|_{H^{-1}}.$$

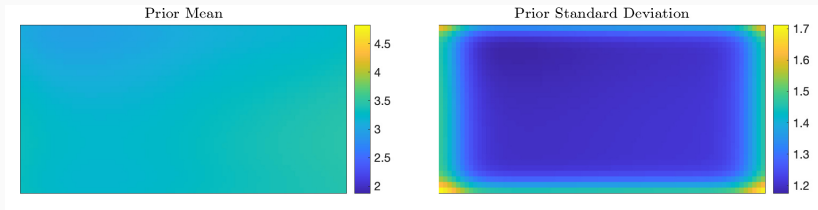
$$d_H(\pi_{\Phi_{L^2}}^y, \pi_{\Phi_{L^2}}^{y'}) \leq C_{L^2} \|y - y'\|_{L^2}.$$

If  $y - y' \approx \sin(kx)$ ,  $\|y - y'\|_{H^{-1}} \approx \mathcal{O}(\frac{1}{k})$ , while  $\|y - y'\|_{L^2} \approx \mathcal{O}(1)$ .

## $W_2$ Likelihood Function — Example

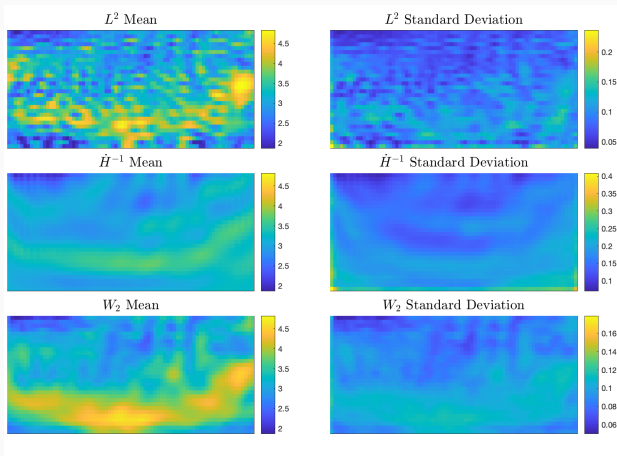


The true continuous velocity field  $v$  and the state parameter  $u = F^{-1}(1/v^2)$ .



[Dunlop-Y,2020] The prior mean  $m_o(x)$  and standard deviation.

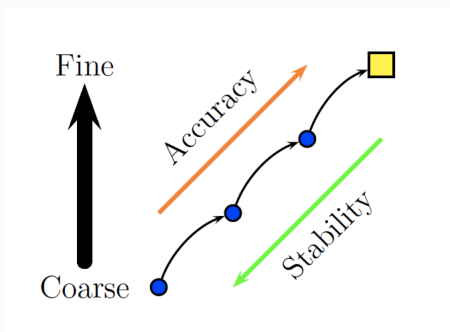
## $W_2$ Likelihood Function — Example



The means (left) and standard deviations (right) of the Laplace approximations.

# Properties of the Wasserstein Metric in Inverse Data Matching

1. Better convexity (optimization landscape) as an objective function for certain problems.
2. Robust with respect to high-frequency noise.
3. As a likelihood function in Bayesian inference for better stability. (Well-posedness of the posterior is proved.)



From [Qiu, 2013]

## Acknowledgment

All my collaborators.



NYU | COURANT

Mathematics



Thank you for the attention!