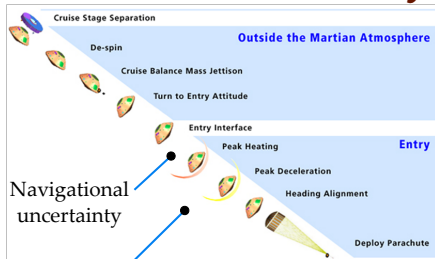# Filtering as Gradient Descent

## Abhishek Halder

Department of Applied Mathematics and Statistics
University of California, Santa Cruz
Santa Cruz, CA 95064

# Joint work with Tryphon T. Georgiou

# Motivation: Mars Entry-Descent-Landing



Image credit: NASA JPL

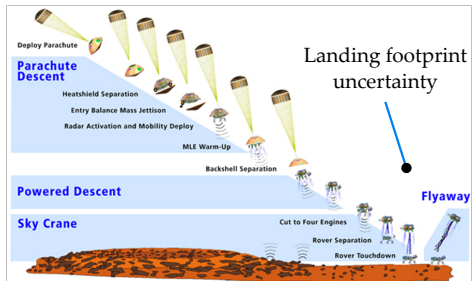# Motivation: Mars Entry-Descent-Landing



Image credit: NASA JPL

Large number of uncertain scenarios ⇝ Probability density

# Motivation: Mars Entry-Descent-Landing



**Supersonic parachute**



**Gale Crater (4.49S, 137.42E)**

# Problem: Uncertainty Propagation



Initial conditions

Process noise

Parameters

Process model

State density

$\rho(\mathbf{x}(t), t)$

# Problem: Uncertainty Propagation



**Trajectory flow:**

$$d\mathbf{X}(t) = \mathbf{f}(\mathbf{X}, t)\, dt + \mathbf{g}(\mathbf{X}, t)\, d\mathbf{w}(t), \quad d\mathbf{w}(t) \sim \mathcal{N}(0, \mathbf{Q}dt)$$

# Problem: Uncertainty Propagation



**Trajectory flow:**

$$d\mathbf{X}(t) = \mathbf{f}(\mathbf{X}, t)\, dt + \mathbf{g}(\mathbf{X}, t)\, d\mathbf{w}(t), \quad d\mathbf{w}(t) \sim \mathcal{N}(0, \mathbf{Q}dt)$$

**Density flow:**

$$\frac{\partial \rho}{\partial t} = \mathcal{L}_{\mathrm{FP}}(\rho) := -\nabla \cdot (\rho \mathbf{f}) + \frac{1}{2} \sum_{i,j=1}^{n} \frac{\partial^2}{\partial x_i \partial x_j} \left( \left( \mathbf{g}\mathbf{Q}\mathbf{g}^\top \right)_{ij} \rho \right)$$

# Problem: Filtering



Initial conditions

Parameters

Process noise

Process model

Prior

Sensor noise

Measurement model

Posterior density

$\rho^+ \left( \mathbf{x}(t), t \right)$

Init condit

Param

# Problem: Filtering



**Trajectory flow:**

$$d\mathbf{x}(t) = \mathbf{f}(\mathbf{x}, t)\, dt + \mathbf{g}(\mathbf{x}, t)\, d\mathbf{w}(t), \quad d\mathbf{w}(t) \sim \mathcal{N}(0, \mathbf{Q}dt)$$
$$d\mathbf{z}(t) = \mathbf{h}(\mathbf{x}, t)\, dt + d\mathbf{v}(t), \qquad\quad d\mathbf{v}(t) \sim \mathcal{N}(0, \mathbf{R}dt)$$
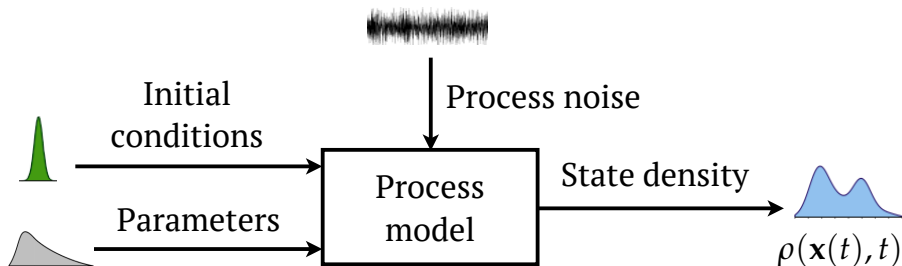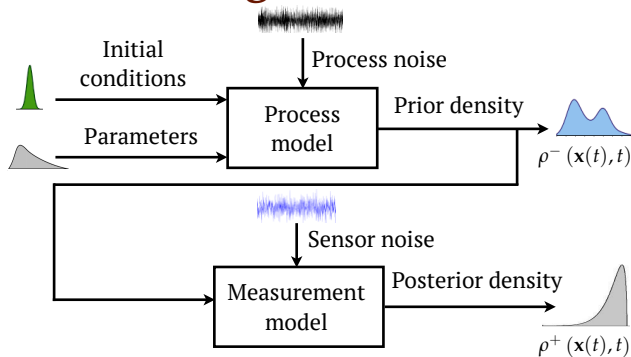
# Problem: Filtering



**Trajectory flow:**

$$d\mathbf{x}(t) = \mathbf{f}(\mathbf{X}, t)\, dt + \mathbf{g}(\mathbf{X}, t)\, d\mathbf{w}(t), \quad d\mathbf{w}(t) \sim \mathcal{N}(0, \mathbf{Q}dt)$$
$$d\mathbf{z}(t) = \mathbf{h}(\mathbf{X}, t)\, dt + d\mathbf{v}(t), \qquad\qquad d\mathbf{v}(t) \sim \mathcal{N}(0, \mathbf{R}dt)$$

**Density flow:**

$$d\rho^+ = \left[ \mathcal{L}_{\mathrm{FP}}dt + \left(\mathbf{h}(\mathbf{x}, t) - \mathbb{E}_{\rho^+}\{\mathbf{h}(\mathbf{x}, t)\}\right)^{\top} \mathbf{R}^{-1}\left(d\mathbf{z}(t) - \mathbb{E}_{\rho^+}\{\mathbf{h}(\mathbf{x}, t)\}dt\right) \right]\rho^+$$

# Research Scope

Density Flow

PDE formulation $\Longleftrightarrow$ Variational formulation

Numerically approximate
solution of PDE

Recursively evaluate
proximal operators

# Research Scope

Density Flow

PDE formulation $\iff$ Variational formulation

Numerically approximate
solution of PDE

Recursively evaluate
proximal operators

Density flow $\rightsquigarrow$ gradient descent in infinite dimensions

# Gradient Descent in Finite Dimensions

**Problem:** $\displaystyle\minimize_{\mathbf{x}\in\mathbb{R}^n}\ \ \phi(\mathbf{x})$

**Algorithm:** $\mathbf{x}_k = \mathbf{x}_{k-1} - h\nabla\phi(\mathbf{x}_{k-1})$

# Gradient Descent in Finite Dimensions

**Problem:** $\displaystyle\operatorname*{minimize}_{\mathbf{x}\in\mathbb{R}^n} \phi(\mathbf{x})$

**Algorithm:** $\mathbf{x}_k = \mathbf{x}_{k-1} - h\nabla\phi(\mathbf{x}_{k-1})$

**Advantage:**

- is a descent method: $\phi(\mathbf{x}_k) \leq \phi(\mathbf{x}_{k-1})$

- convergence under very few assumptions

- simple first order method

- can account constraints (projected gradient descent)

# Why does gradient descent work?

$$z = \phi(x), \quad x \in \mathbb{R}^2$$



$-\nabla\phi(\mathbf{x})$ is the max-rate descending direction (why?)

# Rate of Convergence for Gradient Descent

| If | then |
|---|---|
| $\phi$ is $(\frac{1}{h})$-smooth  ($\Leftrightarrow \nabla\phi$ is $\frac{1}{h}$ Lipschitz) | $O(\frac{1}{kh})$ |
| $\phi$ is $(\frac{1}{h})$-smooth  AND $\sigma$-strongly convex | $O(\frac{1}{h}\exp(-\frac{h\sigma}{2}k))$ |

# Gradient Descent ⟿ Gradient Flow

- GD is **Euler discretization** of GF

$$\frac{d\mathbf{x}}{dt} = -\nabla\phi(\mathbf{x}), \ \mathbf{x} \in \mathbb{R}^n$$

- **Rate matching:**

  GD rate $O(\frac{1}{kh})$ when $\phi$ is $(\frac{1}{h})$-smooth

  GF rate $O(\frac{1}{t})$ when $\phi$ is convex

# Gradient Descent ⤳ Proximal Operator

$$\mathbf{x}_k = \mathbf{x}_{k-1} - h\nabla\phi(\mathbf{x}_{k-1})$$

$$\Updownarrow$$

$$\mathbf{x}_k = \text{proximal}_{h\phi}^{\|\cdot\|}(\mathbf{x}_{k-1})$$

$$:= \underset{\mathbf{x}\in\mathbb{R}^n}{\arg\min} \left\{ \tfrac{1}{2}\|\mathbf{x}-\mathbf{x}_{k-1}\|^2 + h\phi(\mathbf{x}) \right\}$$

# Gradient Descent ⤳ Proximal Operator

$$\mathbf{x}_k = \mathbf{x}_{k-1} - h\nabla\phi(\mathbf{x}_{k-1})$$
$$\Updownarrow$$
$$\mathbf{x}_k = \text{proximal}_{h\phi}^{\|\cdot\|}(\mathbf{x}_{k-1})$$

$$:= \underset{\mathbf{x}\in\mathbb{R}^n}{\arg\min}\left\{\tfrac{1}{2}\|\mathbf{x} - \mathbf{x}_{k-1}\|^2 + h\phi(\mathbf{x})\right\}$$

**This is nice because**

- argmin of $\phi \equiv$ fixed point of prox. operator

- prox. is smooth even when $\phi$ is not

-
> reveals metric structure of gradient descent

# Gradient Descent in Infinite Dimensions



$$z = \Phi(\rho), \quad \rho \in \mathscr{D}$$

$d(\rho_0, \rho_1)$

$\rho_4 \rho_3 \rho_2 \quad \rho_1 \quad \rho_0$

**Proximal recursion:** $\rho_k = \underset{\rho \in \mathscr{D}}{\arginf} \left\{ \frac{1}{2} d^2(\rho, \rho_{k-1}) + h\Phi(\rho) \right\}$

# Gradient Descent Summary

**Finite dimensions**

$$\frac{d\mathbf{x}}{dt} = -\nabla\phi(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^n$$

$\mathbf{x}_k(h) = \mathbf{x}_{k-1} - h\nabla\phi(\mathbf{x}_{k-1})$

$= \underset{\mathbf{x}}{\text{argmin}}\{\frac{1}{2}\|\mathbf{x} - \mathbf{x}_{k-1}\|^2 + h\phi(\mathbf{x})\}$

$= \text{proximal}_{h\phi}^{\|\cdot\|}(\mathbf{x}_{k-1})$

$$\mathbf{x}_k(h) \to \mathbf{x}(t=kh), \text{ as } h \downarrow 0$$

**Infinite dimensions**

$$\frac{\partial\rho}{\partial t} = \mathcal{L}(\mathbf{x}, \rho), \quad \mathbf{x} \in \mathbb{R}^n, \rho \in \mathscr{D}$$

$\rho_k(\mathbf{x}, h)$

$= \underset{\rho}{\text{argmin}}\{\frac{1}{2}d(\rho, \rho_{k-1})^2 + h\Phi(\rho)\}$

$= \text{proximal}_{h\Phi}^{d(\cdot,\cdot)}(\rho_{k-1})$

$$\rho_k(\mathbf{x}, h) \rightharpoonup \rho(\mathbf{x}, t=kh), \text{ as } h \downarrow 0$$

# Related Work

| Transport PDE $\frac{\partial \rho}{\partial t} = \mathcal{L}(\mathbf{x}, \rho)$ | Gradient descent scheme | |
|:---:|:---:|:---:|
| $\mathcal{L}(\mathbf{x}, \rho)$ | $\frac{1}{2} d^2(\rho, \rho_{k-1})$ | $\Phi(\rho)$ |
| $\triangle \rho$ <br><br> Heat equation (1822) | $\frac{1}{2} \parallel \rho - \rho_{k-1} \parallel_{L_2(\mathbb{R}^n)}^2$ <br><br> Squared $L_2$ norm of difference | $\frac{1}{2} \int_{\mathbb{R}^n} \parallel \nabla \rho \parallel^2$ <br><br> Dirichlet energy, CFL (1928) |
| $\nabla \cdot (\nabla U(\mathbf{x}) \rho) + \beta^{-1} \triangle \rho$ <br><br> Fokker-Planck-Kolmogorov PDE (1914,'17,'31) | $\frac{1}{2} W^2(\rho, \rho_{k-1})$ <br><br> Optimal transport cost | $\mathbb{E}_\rho \left[ U(\mathbf{x}) + \beta^{-1} \log \rho \right]$ <br><br> Free energy, JKO (1998) |
| $\left( (\mathbf{h} - \mathbb{E}_\rho[\mathbf{h}])^\top \mathbf{R}^{-1} (\mathrm{d}\mathbf{z} - \mathbb{E}_\rho[\mathbf{h}]\mathrm{d}t) \right) \rho$ <br><br> Kushner-Stratonovich SPDE (1964,'59) | $\mathrm{D}_{KL}(\rho || \rho_{k-1})$ <br><br> Kullback-Leibler divergence | $\frac{1}{2} \mathbb{E}_\rho[(\mathbf{y}_k - \mathbf{h})^\top \mathbf{R}^{-1} (\mathbf{y}_k - \mathbf{h})]$ <br><br> Quadratic surprise, LMMR (2015) |

# Related Work

| Transport PDE $\frac{\partial \rho}{\partial t} = \mathcal{L}(\mathbf{x}, \rho)$ | Gradient descent scheme | |
|---|---|---|
| $\mathcal{L}(\mathbf{x}, \rho)$ | $\frac{1}{2} d^2(\rho, \rho_{k-1})$ | $\Phi(\rho)$ |
| $\triangle \rho$ <br><br> Heat equation (1822) | $\frac{1}{2} \\| \rho - \rho_{k-1} \\|^2_{L_2(\mathbb{R}^n)}$ <br><br> Squared $L_2$ norm of difference | $\frac{1}{2} \int_{\mathbb{R}^n} \\| \nabla \rho \\|^2$ <br><br> Dirichlet energy, CFL (1928) |
| $\nabla \cdot (\nabla U(\mathbf{x})\rho) + \beta^{-1} \triangle \rho$ <br><br> Fokker-Planck-Kolmogorov PDE (1914,'17,'31) | $\frac{1}{2} W^2(\rho, \rho_{k-1})$ <br><br> Optimal transport cost | $\mathbb{E}_\rho \left[ U(\mathbf{x}) + \beta^{-1} \log \rho \right]$ <br><br> Free energy, JKO (1998) |
| $\left( (\mathbf{h} - \mathbb{E}_\rho[\mathbf{h}])^\top \mathbf{R}^{-1} (d\mathbf{z} - \mathbb{E}_\rho[\mathbf{h}]dt) \right) \rho$ <br><br> Kushner-Stratonovich SPDE (1964,'59) | $D_{KL}(\rho \|| \rho_{k-1})$ <br><br> Kullback-Leibler divergence | $\frac{1}{2} \mathbb{E}_\rho[(\mathbf{y}_k - \mathbf{h})^\top \mathbf{R}^{-1}(\mathbf{y}_k - \mathbf{h})]$ <br><br> Quadratic surprise, LMMR (2015) |

Process dynamics is stochastic gradient flow:

Gibbs density

$$d\mathbf{x}(t) = -\nabla U(\mathbf{x}) \, dt + \sqrt{2\beta^{-1}} d\mathbf{w}(t), \qquad \rho_\infty(\mathbf{x}) \propto e^{-\beta U(\mathbf{x})}$$

# Related Work

| Transport PDE $\frac{\partial \rho}{\partial t} = \mathcal{L}(\mathbf{x}, \rho)$ | Gradient descent scheme | |
|---|---|---|
| $\mathcal{L}(\mathbf{x}, \rho)$ | $\frac{1}{2} d^2(\rho, \rho_{k-1})$ | $\Phi(\rho)$ |
| $\triangle \rho$ <br><br> Heat equation (1822) | $\frac{1}{2} \parallel \rho - \rho_{k-1} \parallel^2_{L_2(\mathbb{R}^n)}$ <br><br> Squared $L_2$ norm of difference | $\frac{1}{2} \int_{\mathbb{R}^n} \parallel \nabla \rho \parallel^2$ <br><br> Dirichlet energy, CFL (1928) |
| $\nabla \cdot (\nabla U(\mathbf{x})\rho) + \beta^{-1} \triangle \rho$ <br><br> Fokker-Planck-Kolmogorov PDE (1914,'17,'31) | $\frac{1}{2} W^2(\rho, \rho_{k-1})$ <br><br> Optimal transport cost | $\mathbb{E}_\rho \left[ U(\mathbf{x}) + \beta^{-1} \log \rho \right]$ <br><br> Free energy, JKO (1998) |
| $\left( (\mathbf{h} - \mathbb{E}_\rho[\mathbf{h}])^\top \mathbf{R}^{-1} (d\mathbf{z} - \mathbb{E}_\rho[\mathbf{h}]dt) \right) \rho$ <br><br> Kushner-Stratonovich SPDE (1964,'59) | $D_{KL}(\rho || \rho_{k-1})$ <br><br> Kullback-Leibler divergence | $\frac{1}{2} \mathbb{E}_\rho [(\mathbf{y}_k - \mathbf{h})^\top \mathbf{R}^{-1} (\mathbf{y}_k - \mathbf{h})]$ <br><br> Quadratic surprise, LMMR (2015) |

No process dynamics, only measurement update:

$$d\mathbf{x}(t) = 0, \quad d\mathbf{z}(t) = \mathbf{h}(\mathbf{x}, t) \, dt + d\mathbf{v}(t), \quad d\mathbf{v}(t) \sim \mathcal{N}(0, \mathbf{R}dt)$$

## Our Contribution

| Transport description | Gradient descent scheme | |
|---|---|---|
| PDE/SDE/ODE | $\frac{1}{2}d^2(\rho, \rho_{k-1})$ | $\Phi(\rho)$ |
| Mean ODE, Lyapunov ODE <br><br> Linear Gaussian uncertainty propagation | $\frac{1}{2}W^2(\rho, \rho_{k-1})$ <br><br> Optimal transport cost | $\mathbb{E}_\rho\left[U(\mathbf{x}, t) + \frac{\text{tr}(\mathbf{P}_\infty)}{n}\log\rho\right]$ <br><br> Generalized free energy |
| Conditional mean SDE, Riccati ODE <br><br> Kalman-Bucy filter | $D_{KL}(\rho\|\|\rho_{k-1})$ <br><br> Kullback-Leibler divergence | $\frac{1}{2}\mathbb{E}_\rho[(\mathbf{y}_k - \mathbf{h})^\top \mathbf{R}^{-1}(\mathbf{y}_k - \mathbf{h})]$ <br><br> Quadratic surprise |
| ditto | $\frac{1}{2}d_{\text{FR}}^2(\rho, \rho_{k-1})$ <br><br> Fisher-Rao metric | ditto |
| Kushner-Stratonovich SPDE <br><br> Nonlinear filter | ditto <br><br> Fisher-Rao metric | ditto |

# The Case for Linear Gaussian Systems

**Model:**

$$d\mathbf{x}(t) = \mathbf{A}\mathbf{x}(t)dt + \mathbf{B}d\mathbf{w}(t), \quad d\mathbf{w}(t) \sim \mathcal{N}(0, \mathbf{Q}dt)$$

$$d\mathbf{z}(t) = \mathbf{C}\mathbf{x}(t)dt + d\mathbf{v}(t), \qquad d\mathbf{v}(t) \sim \mathcal{N}(0, \mathbf{R}dt)$$

**Given $\mathbf{x}(0) \sim \mathcal{N}(\mu_0, \mathbf{P}_0)$, want to recover:**

For uncertainty propagation:

$$\dot{\mu} = \mathbf{A}\mu, \ \mu(0) = \mu_0; \quad \dot{\mathbf{P}} = \mathbf{A}\mathbf{P} + \mathbf{P}\mathbf{A}^\top + \mathbf{B}\mathbf{Q}\mathbf{B}^\top, \ \mathbf{P}(0) = \mathbf{P}_0.$$

For filtering:

$$d\mu^+(t) = \mathbf{A}\mu^+(t)dt + \underset{\substack{| \\ \mathbf{K}(t)}}{\boxed{\mathbf{K}(t)}}^{\mathbf{P}^+\mathbf{C}\mathbf{R}^{-1}} (d\mathbf{z}(t) - \mathbf{C}\mu^+(t)dt),$$

$$\dot{\mathbf{P}}^+(t) = \mathbf{A}\mathbf{P}^+(t) + \mathbf{P}^+(t)\mathbf{A}^\top + \mathbf{B}\mathbf{Q}\mathbf{B}^\top - \mathbf{K}(t)\mathbf{R}\mathbf{K}(t)^\top.$$

# The Case for Linear Gaussian Systems

**Challenge 1:**

How to actually perform the infinite dimensional optimization over $\mathscr{D}_2$?

**Challenge 2:**

If and how one can apply the variational schemes for generic linear system with Hurwitz $\mathbf{A}$ and controllable $(\mathbf{A}, \mathbf{B})$?

# Addressing Challenge 1: How to Compute

## Two Step Optimization Strategy

– Notice that the objective is a *sum*:

first
functional

second
functional

$$\operatorname*{arginf}_{\rho \in \mathscr{D}_2}\{\ \frac{1}{2}d(\rho, \rho_{k-1})^2 \ + \ h\Phi(\rho)\ \}$$

– Choose a parametrized subspace of $\mathscr{D}_2$ such that the individual minimizers over that subspace match

– Then optimize over parameters

– $\mathscr{D}_{\mu,\mathbf{P}} \subset \mathscr{D}_2$ works!

# Addressing Challenge 2: Generic $(\mathbf{A}, \sqrt{2}\mathbf{B})$

## Two Successive Coordinate Transformations

---

**#1. Equipartition of energy:**

- Define *thermodynamic temperature* $\theta := \frac{1}{n}\text{tr}(\mathbf{P}_\infty)$, and *inverse temperature* $\beta := \theta^{-1}$

- State vector: $\mathbf{x} \mapsto \mathbf{x}_{ep} := \sqrt{\theta}\mathbf{P}_\infty^{-\frac{1}{2}}\mathbf{x}$

- System matrices:

$$\mathbf{A}, \sqrt{2}\mathbf{B} \mapsto \underset{\mathbf{A}_{ep}}{\underbrace{\mathbf{P}_\infty^{-\frac{1}{2}}\mathbf{A}\mathbf{P}_\infty^{\frac{1}{2}}}}, \sqrt{2\theta}\ \underset{\mathbf{B}_{ep}}{\underbrace{\mathbf{P}_\infty^{-\frac{1}{2}}\mathbf{B}}}$$

- Stationary covariance:
$\mathbf{P}_\infty \mapsto \theta\mathbf{I}$

# Addressing Challenge 2: Generic $(\mathbf{A}, \sqrt{2}\mathbf{B})$

## Two Successive Coordinate Transformations

**#2. Symmetrization:**

– State vector: $\mathbf{x}_{\text{ep}} \mapsto \mathbf{x}_{\text{sym}} := e^{-\mathbf{A}_{\text{ep}}^{\text{skew}}t}\mathbf{x}_{\text{ep}}$

– System matrices:

$$\mathbf{A}_{\text{ep}}, \sqrt{2\theta}\mathbf{B}_{\text{ep}} \mapsto \underbrace{e^{-\mathbf{A}_{\text{ep}}^{\text{skew}}t}\mathbf{A}_{\text{ep}}^{\text{sym}}e^{\mathbf{A}_{\text{ep}}^{\text{skew}}t}}_{\mathbf{F}(t)}, \sqrt{2\theta}\,\underbrace{e^{-\mathbf{A}_{\text{ep}}^{\text{skew}}t}\mathbf{B}_{\text{ep}}}_{\mathbf{G}(t)}$$

– Stationary covariance:
$\theta\mathbf{I} \mapsto \theta\mathbf{I}$

– Potential: $U(\mathbf{x}_{\text{sym}}, t) := -\frac{1}{2}\mathbf{x}_{\text{sym}}^{\top}\mathbf{F}(t)\mathbf{x}_{\text{sym}} \geq 0$

# Summary

– Two successive coordinate transformations bring generic linear system to JKO canonical form

– Can apply two step optimization strategy in $\mathbf{x}_{\text{sym}}$ coordinate

– Recovers mean-covariance propagation, and Kalman-Bucy filter in $h \downarrow 0$ limit

– Changing the distance in LMMR from $D_{\text{KL}}$ to $\frac{1}{2}W_2^2$ gives Luenberger-type observers

– **Future work:** computation for nonlinear filtering

# Thank You

# Backup Slides

# Gradient Descent with Constraints

$$\underset{\mathbf{x} \in \mathcal{C}}{\text{minimize}} \quad \phi(\mathbf{x})$$

$$\Updownarrow$$

$$\mathbf{x}_k = \text{proj}_{\mathcal{C}} \left( \mathbf{x}_{k-1} - h \nabla \phi(\mathbf{x}_{k-1}) \right)$$

$$\Updownarrow$$

$$\mathbf{x}_k = \text{proximal}_{h\phi}^{\|\cdot\|} \left( \mathbf{x}_{k-1} \right)$$

$$:= \underset{\mathbf{x} \in \mathcal{C}}{\text{argmin}} \left\{ \tfrac{1}{2} \|\mathbf{x} - \mathbf{x}_{k-1}\|^2 + h\phi(\mathbf{x}) \right\}$$