

# Deep Learning for Inverse Problems

Where are we? How far can we go?

---

Jonas Adler<sup>1, 2</sup>    Ozan Öktem<sup>1</sup>

<sup>1</sup>Department of Mathematics

KTH - Royal Institute of Technology, Stockholm, Sweden

<sup>2</sup>Research and Physics

Elekta, Stockholm, Sweden



$$y = \mathcal{A}(x^*) + e.$$

$$y \in Y$$

Data

$$x^* \in X$$

Image

$$\mathcal{A} : X \rightarrow Y$$

Forward operator

$$e \in Y$$

Noise

$$y = \mathcal{A}(x^*) + e.$$

$$y \in Y$$

$$x^* \in X$$

$$\mathcal{A} : X \rightarrow Y$$

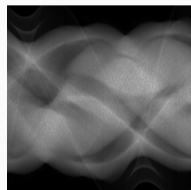
$$e \in Y$$

Data

Image

Forward operator

Noise



$$y = \mathcal{A}(x^*) + e.$$

$y \in Y$

$x^* \in X$

$\mathcal{A} : X \rightarrow Y$

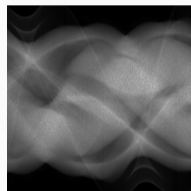
$e \in Y$

Data

Image

Forward operator

Noise



$$y = \mathcal{A}(x^*) + e.$$

$$y \in Y$$

$$x^* \in X$$

$$\mathcal{A}: X \rightarrow Y$$

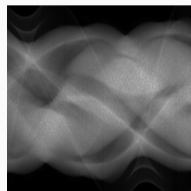
$$e \in Y$$

Data

Image

Forward operator

Noise



$$y = \mathcal{A}(x^*) + e.$$

$$y \in Y$$

$$x^* \in X$$

$$\mathcal{A} : X \rightarrow Y$$

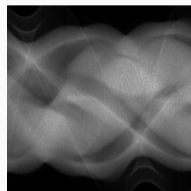
$$e \in Y$$

Data

Image

Forward operator

Noise



$$y = \mathcal{A}(x^*) + e.$$

$$y \in Y$$

$$x^* \in X$$

$$\mathcal{A} : X \rightarrow Y$$

$$e \in Y$$

Data

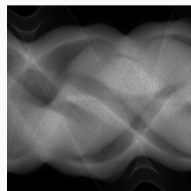
Image

Forward operator

Noise



$$\begin{array}{c} \xrightarrow{\mathcal{A}} \\ \xleftarrow{\text{"}\mathcal{A}^{-1}\text{"}} \end{array}$$



$$y = \mathcal{A}(x^*) + e.$$

$$y \in Y$$

$$x^* \in X$$

$$\mathcal{A} : X \rightarrow Y$$

$$e \in Y$$

Data

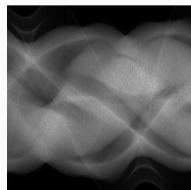
Image

Forward operator

Noise



$$\begin{array}{c} \xrightarrow{\mathcal{A}} \\ \xleftarrow{\mathcal{A}^{-1}} \end{array}$$



The problem is ill-posed: non-uniqueness, instability

Data  $y \in Y$  is a single observation generated by  $Y$ -valued random variable  $\mathbf{y}$  where

$$\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}.$$

Data  $y \in Y$  is a single observation generated by  $Y$ -valued random variable  $\mathbf{y}$  where

$$\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}.$$

**Full solution:** A probability distribution on model parameter space  $X$

$$\mathbb{P}(\mathbf{x} \mid \mathbf{y} = y)$$

This is a full characterization of the reconstruction, including uncertainty.

Data  $y \in Y$  is a single observation generated by  $Y$ -valued random variable  $\mathbf{y}$  where

$$\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}.$$

**Full solution:** A probability distribution on model parameter space  $X$

$$\mathbb{P}(\mathbf{x} \mid \mathbf{y} = y)$$

This is a full characterization of the reconstruction, including uncertainty.

**Typical solution:** Compute some estimator, e.g. the conditional mean

$$\mathbb{E}[\mathbf{x} \mid \mathbf{y} = y]$$

## Theorem (Conditional Mean)

*Assume that  $Y$  is a measurable metric space,  $X$  a measurable Hilbert space, and  $\mathbf{y}$  and  $\mathbf{x}$  are  $Y$ - and  $X$ -valued random variables, respectively. Let*

$$h^* = \arg \min_{h: Y \rightarrow X} \mathbb{E} \left[ \|h(\mathbf{y}) - \mathbf{x}\|_X^2 \right].$$

*Then  $h^*(y) := \mathbb{E}[\mathbf{x} \mid \mathbf{y} = y]$  almost everywhere.*

- Suppose we aim to compute

$$\mathbb{E}[\mathbf{x} \mid \mathbf{y} = y]$$

- This can be done by solving

- Suppose we aim to compute

$$\mathbb{E}[\mathbf{x} \mid \mathbf{y} = y]$$

- This can be done by solving

$$\min_{h: Y \rightarrow X} \mathbb{E} \left[ \|h(\mathbf{y}) - \mathbf{x}\|_X^2 \right].$$

- Suppose we aim to compute

$$\mathbb{E}[\mathbf{x} \mid \mathbf{y} = y]$$

- This can be done by solving

$$\min_{h: Y \rightarrow X} \mathbb{E} \left[ \|h(\mathbf{y}) - \mathbf{x}\|_X^2 \right].$$

The minimization is over all measurable functions

- Suppose we aim to compute

$$\mathbb{E}[\mathbf{x} \mid \mathbf{y} = y]$$

- This can be done by solving

$$\min_{\theta \in \Theta} \mathbb{E} \left[ \|\mathcal{A}_{\theta}^{\dagger}(\mathbf{y}) - \mathbf{x}\|_X^2 \right].$$

The minimization is over all measurable functions

Restrict minimization to some tractable subset

- Suppose we aim to compute

$$\mathbb{E}[\mathbf{x} \mid \mathbf{y} = y]$$

- This can be done by solving

$$\min_{\theta \in \Theta} \mathbb{E} \left[ \|\mathcal{A}_{\theta}^{\dagger}(\mathbf{y}) - \mathbf{x}\|_X^2 \right].$$

Expectation is taken over the unknown joint distribution

- Suppose we aim to compute

$$\mathbb{E}[\mathbf{x} \mid \mathbf{y} = y]$$

- This can be done by solving

$$\min_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N \|\mathcal{A}_{\theta}^{\dagger}(y_i) - x_i\|_X^2$$

Expectation is taken over the unknown joint distribution

Replace with empirical mean

- Suppose we aim to compute

$$\mathbb{E}[\mathbf{x} \mid \mathbf{y} = y]$$

- This can be done by solving

$$\min_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N \|\mathcal{A}_{\theta}^{\dagger}(y_i) - x_i\|_X^2$$

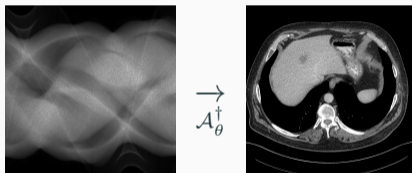
This is a "computationally tractable" formulation, we just need to pick  $\{\mathcal{A}_{\theta}^{\dagger}\}_{\theta \in \Theta}$ .

Architecture: Specification of the class of operators  $\{\mathcal{A}_\theta^\dagger\}_{\theta \in \Theta}$ .

# Learned inversion methods

Architecture: Specification of the class of operators  $\{\mathcal{A}_\theta^\dagger\}_{\theta \in \Theta}$ .

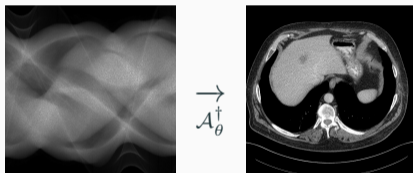
Main complication:  $\mathcal{A}_\theta^\dagger : Y \rightarrow X$ .



# Learned inversion methods

Architecture: Specification of the class of operators  $\{\mathcal{A}_\theta^\dagger\}_{\theta \in \Theta}$ .

Main complication:  $\mathcal{A}_\theta^\dagger : Y \rightarrow X$ .

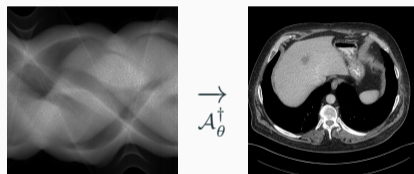


- **Fully learned:** Learn everything, disregard structure.

# Learned inversion methods

Architecture: Specification of the class of operators  $\{\mathcal{A}_\theta^\dagger\}_{\theta \in \Theta}$ .

Main complication:  $\mathcal{A}_\theta^\dagger : Y \rightarrow X$ .

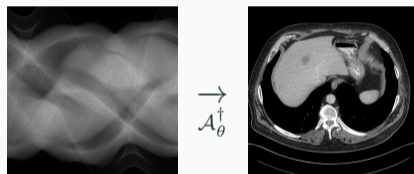


- Fully learned: Learn everything, disregard structure.
- **Learned post-processing**: First apply standard inverse, then denoise  $\mathcal{A}_\theta^\dagger = P_\theta \circ \mathcal{A}^\dagger$

# Learned inversion methods

Architecture: Specification of the class of operators  $\{\mathcal{A}_\theta^\dagger\}_{\theta \in \Theta}$ .

Main complication:  $\mathcal{A}_\theta^\dagger : Y \rightarrow X$ .



- Fully learned: Learn everything, disregard structure.
- Learned post-processing: First apply standard inverse, then denoise  $\mathcal{A}_\theta^\dagger = P_\theta \circ \mathcal{A}^\dagger$
- **Learned iterative schemes**: Embed physics inside deep neural network

How well does this actually work?

Measure *generalization gap*:

$$\mathbb{E} \left[ \left\| \mathcal{A}_{\theta^*}^\dagger(\mathbf{y}) - \mathbf{x} \right\|_X^2 \right] - \mathbb{E} \left[ \left\| \mathbb{E}[\mathbf{x} \mid \mathbf{y}] - \mathbf{x} \right\|_X^2 \right].$$

## Results on toy case

Results for ray transform inversion in 2D:

- Inverse problem:

$$y = \mathcal{A}(x) + e$$

- Geometry: Parallel beam, sparse view (30 angles)
- Noise: 5% additive Gaussian
- Training data:  $128 \times 128$  pixel ellipses

## Results on toy case

Results for ray transform inversion in 2D:

- Inverse problem:

$$y = \mathcal{A}(x) + e$$

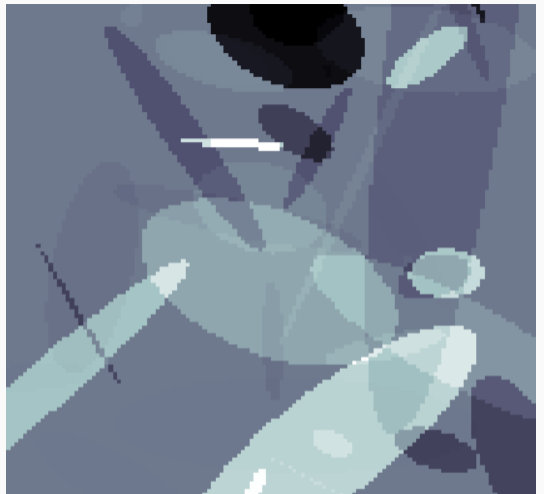
- Geometry: Parallel beam, sparse view (30 angles)
- Noise: 5% additive Gaussian
- Training data:  $128 \times 128$  pixel ellipses

Compare to:

- FBP
- Total Variation
- Post-processing deep learning by U-Net
- Conditional expectation,  $\mathbb{E}(x | y)$ , via MCMC

Measure relative error:

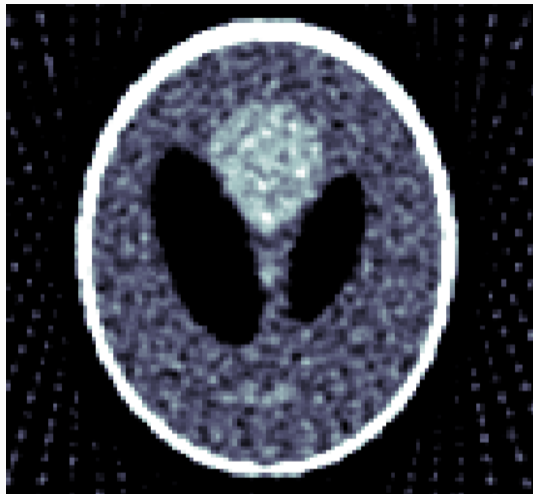
$$\frac{\mathbb{E} \left[ \left\| \mathcal{A}_{\theta^*}^\dagger(y) - x \right\|_X^2 \right]}{\mathbb{E} \left[ \left\| \mathbb{E}[x | y] - x \right\|_X^2 \right]}$$



Training data



Phantom



FBP

Normalized error: 372



Phantom



TV

Normalized error: 56.0



Phantom



Learned Post-processing  
Normalized error: 42.2



Phantom



Learned Iterative  
Normalized error: 5.2



Phantom



Conditional Expectation

Normalized error: 1

# Conclusions on learned reconstruction

- We can find a reconstruction operator by solving a minimization problem
- Architecture: Specification of the class of operators  $\{\mathcal{A}_\theta^\dagger\}_{\theta \in \Theta}$ .
- Learning:

$$\min_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N \|\mathcal{A}_\theta^\dagger(y_i) - x_i\|_X^2$$

- Empirically, current methods are remarkably close to optimal

# What now?

- Apparently deep learning techniques are great for the conditional mean
- What about other estimators?
- Maximum a-posteriori is very hard

# What now?

- Apparently deep learning techniques are great for the conditional mean
- What about other estimators?
- Maximum a-posteriori is very hard
- But, what about finding the whole posterior?

# Generative Adversarial Networks

- Main idea: train two networks, generator  $G$  and discriminator  $D$
- Generator tries to generate "true" samples, discriminator tries to say "good/bad"

# Generative Adversarial Networks

- Main idea: train two networks, generator  $G$  and discriminator  $D$
- Generator tries to generate "true" samples, discriminator tries to say "good/bad"



# Conditional Wasserstein GAN

Input: Supervised training data  $(x_i, y_i)$  generated by  $(\mathbf{x}, \mathbf{y})$ .

Goal: Sample from unknown posterior  $\mathbb{P}(\mathbf{x} \mid \mathbf{y} = y)$ .

Approach: Learn how to sample from posterior by solving

$$\min_{\theta} \mathbb{E}_{\mathbf{y} \sim \mathbb{P}_{data}} \left[ \mathcal{W}(G_{\theta}(\mathbf{y}), \mathbb{P}(\mathbf{x} \mid \mathbf{y})) \right].$$

We minimize the *Wasserstein* distance between the random variables  $G_{\theta}(\mathbf{y})$  and  $\mathbb{P}(\mathbf{x} \mid \mathbf{y})$ !

# Conditional Wasserstein GAN

Input: Supervised training data  $(x_i, y_i)$  generated by  $(\mathbf{x}, \mathbf{y})$ .

Goal: Sample from unknown posterior  $\mathbb{P}(\mathbf{x} \mid \mathbf{y} = y)$ .

Approach: Learn how to sample from posterior by solving

$$\min_{\theta} \mathbb{E}_{\mathbf{y} \sim \mathbb{P}_{data}} \left[ \mathcal{W}(G_{\theta}(\mathbf{y}), \mathbb{P}(\mathbf{x} \mid \mathbf{y})) \right].$$

Re-write using the Kantorovich-Rubinstein dual characterization of  $\mathcal{W}$ .

# Conditional Wasserstein GAN

Input: Supervised training data  $(x_i, y_i)$  generated by  $(\mathbf{x}, \mathbf{y})$ .

Goal: Sample from unknown posterior  $\mathbb{P}(\mathbf{x} \mid \mathbf{y} = y)$ .

Approach: Learn how to sample from posterior by solving

$$\min_{\theta} \left\{ \max_{D \in Lip(X)} \mathbb{E} \left[ D(\mathbf{x}, \mathbf{y}) - D(G_{\theta}(\mathbf{y})) \right] \right\}.$$

Re-write using the Kantorovich-Rubinstein dual characterization of  $\mathcal{W}$ .

# Conditional Wasserstein GAN

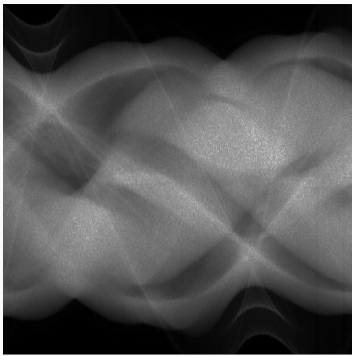
Input: Supervised training data  $(x_i, y_i)$  generated by  $(\mathbf{x}, \mathbf{y})$ .

Goal: Sample from unknown posterior  $\mathbb{P}(\mathbf{x} \mid \mathbf{y} = y)$ .

Approach: Learn how to sample from posterior by solving

$$\min_{\theta} \left\{ \max_{D \in Lip(X)} \mathbb{E} \left[ D(\mathbf{x}, \mathbf{y}) - D(G_{\theta}(\mathbf{y})) \right] \right\}.$$

Formulation useful for deep learning

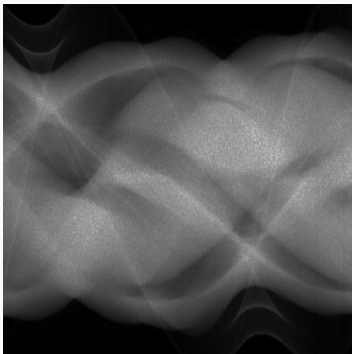


Data



FBP

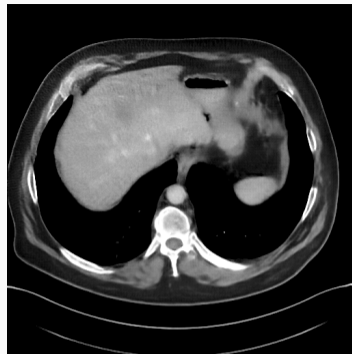
- Case: Patient with suspected metastasis to the liver.
- Data: Clinical helical 3D CT data, 2% of a normal dose (ultra low-dose).



Data

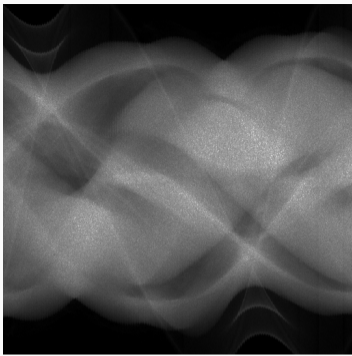


FBP



Posterior mean

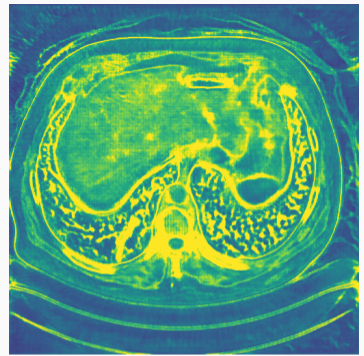
- Case: Patient with suspected metastasis to the liver.
- Data: Clinical helical 3D CT data, 2% of a normal dose (ultra low-dose).



Data

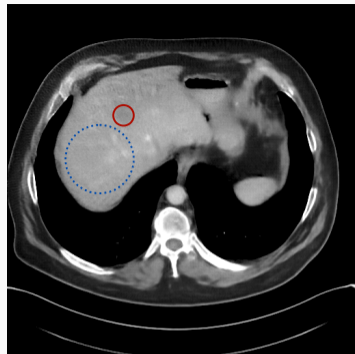


FBP



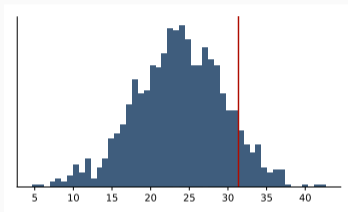
Standard deviation

- Case: Patient with suspected metastasis to the liver.
- Data: Clinical helical 3D CT data, 2% of a normal dose (ultra low-dose).

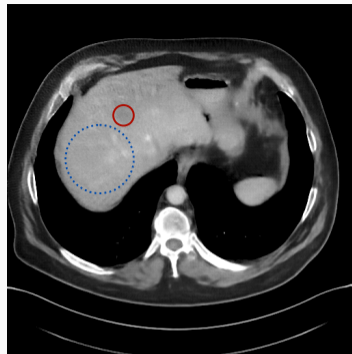


Posterior mean with ROI

- Case: Patient with suspected metastasis to the liver.
- Data: Clinical helical 3D CT data, 2% of a normal dose (ultra low-dose).
- Liver lesion:  $\triangle$  = difference in average contrast between ROI and liver.
- Hypothesis test: Based on 1 000 samples, the ROI contains a lesion at 95%

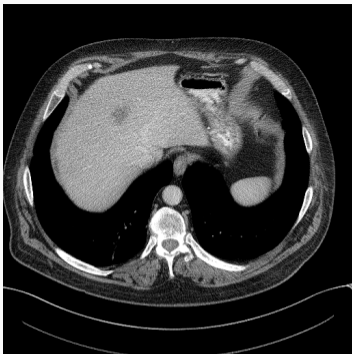


Histogram of  $\Delta$

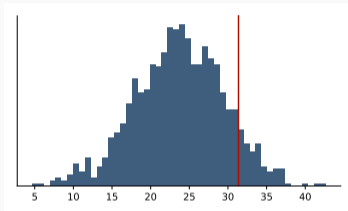


Posterior mean with ROI

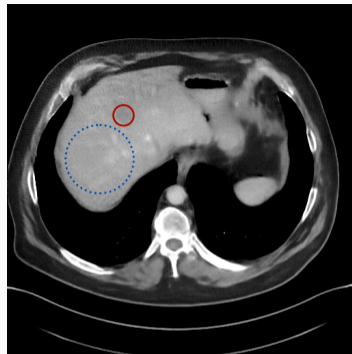
- Case: Patient with suspected metastasis to the liver.
- Data: Clinical helical 3D CT data, 2% of a normal dose (ultra low-dose).
- Liver lesion:  $\Delta$  = difference in average contrast between ROI and liver.
- Hypothesis test: Based on 1 000 samples, the ROI contains a lesion at 95%



Normal dose image



Histogram of  $\Delta$



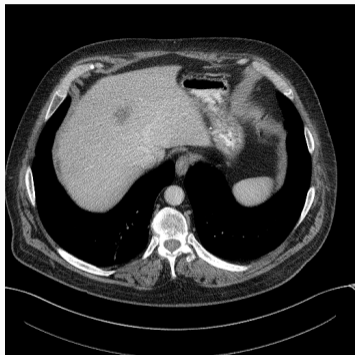
Posterior mean with ROI

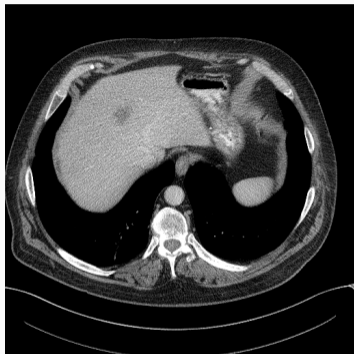
- Case: Patient with suspected metastasis to the liver.
- Data: Clinical helical 3D CT data, 2% of a normal dose (ultra low-dose).
- Liver lesion:  $\Delta$  = difference in average contrast between ROI and liver.
- Hypothesis test: Based on 1 000 samples, the ROI contains a lesion at 95%

- Deep Learning methods for inverse problems building on empirical risk minimization are very powerful
- Fruitful ways forward involve questioning what we're trying to compute
- Posterior sampling is one such option

# Postdoc in Deep Learning based Reconstruction for Spectral-CT

- Theory and methods for machine learning in image reconstruction.
- We've got the worlds first clinical photon counting spectral-CT data.
- Very nice position (great group, travel, salary)
- Pursued jointly with MedTechLabs and the Medical Imaging group at KTH.





Thank you for your attention!