

Deep learning aided Bayesian inference

Deep Ray

Department of Aerospace & Mechanical Engineering
University of Southern California

Email: deepray@usc.edu
Website: deepray.github.io

REU Exposure Seminar
University of Maryland, July 27, 2022

- ▶ Probability, Monte Carlo and Bayes theorem
- ▶ Forward, inverse problems and the Bayesian framework
- ▶ Neural networks
- ▶ Conditional generative adversarial networks (cGANs)
- ▶ Applications:
 - ▶ Academic example: heat equation
 - ▶ Real life example: MRI brain extraction

- ▶ X is an \mathbb{R} -valued random variable with distribution P_X , e.g. uniform $\mathcal{U}([a, b])$, Gaussian $N(\mu, \sigma^2)$.
- ▶ **Expectations:** For a function $f : \mathbb{R} \rightarrow \mathbb{R}$, the expectation of $f(X)$ is

$$\mathbb{E}_{X \sim P_X} [f(X)] = \int_{\mathbb{R}} f(x) P_X(x) \mathrm{d}x.$$

For $f(x) = x$, we get the mean

$$\mu_X := \mathbb{E}_{X \sim P_X} [X] = \int_{\mathbb{R}} x P_X(x) \mathrm{d}x.$$

For $f(x) = (x - \mu_X)^2$, we get the variance

$$\text{Var}(X) := \mathbb{E}_{X \sim P_X} [(X - \mu_X)^2] = \int_{\mathbb{R}} (x - \mu_X)^2 P_X(x) \mathrm{d}x.$$

The standard deviation is $SD(X) := \sqrt{\text{Var}(X)}$.

- ▶ $\mathbf{X} = (X_1, X_2)$ is an \mathbb{R}^2 -valued (2 dimensional) random variable with the **joint** distribution $P_{\mathbf{X}} = P_{X_1 X_2}$, e.g. Gaussian $N(\boldsymbol{\mu}, \Sigma)$.
- ▶ **Expectations:** For a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$

$$\mathbb{E}_{\mathbf{X} \sim P_{\mathbf{X}}} [f(\mathbf{X})] = \int f(\mathbf{x}) P_{\mathbf{X}}(\mathbf{x}) d\mathbf{x}.$$

For a vector-valued function $\mathbf{f} : \mathbb{R}^2 \rightarrow \mathbb{R}^k$ with $\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), \dots, f_k(\mathbf{x})]$,

$$\mathbb{E}_{\mathbf{X} \sim P_{\mathbf{X}}} [\mathbf{f}(\mathbf{X})] = \left[\mathbb{E}_{\mathbf{X} \sim P_{\mathbf{X}}} [f_1(\mathbf{X})], \dots, \mathbb{E}_{\mathbf{X} \sim P_{\mathbf{X}}} [f_k(\mathbf{X})] \right] \in \mathbb{R}^k.$$

In particular, $\mathbf{f}(\mathbf{x}) = \mathbf{x}$ gives 2-dimensional mean

$$\boldsymbol{\mu}_{\mathbf{X}} := \mathbb{E}_{\mathbf{X} \sim P_{\mathbf{X}}} [\mathbf{X}] = [\mu_1, \mu_2].$$

The covariance matrix $\Sigma_{\mathbf{X}} \in \mathbb{R}^{2 \times 2}$ is given by

$$\Sigma_{\mathbf{X}} := \begin{bmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) \\ \text{Cov}(X_1, X_2) & \text{Var}(X_2) \end{bmatrix}.$$

- **Marginals:** The probability distributions of each component

$$P_{X_1}(x_1) = \int_{\mathbb{R}} P_{\mathbf{X}}(x_1, x_2) dx_2, \quad P_{X_2}(x_2) = \int_{\mathbb{R}} P_{\mathbf{X}}(x_1, x_2) dx_1.$$

- **Conditionals:** The conditional probability distribution of X_1 given $X_2 = x_2$

$$P_{X_1|X_2}(x_1|x_2) = \frac{P_{X_1X_2}(x_1, x_2)}{P_{X_2}(x_2)}.$$

Similarly, the conditional probability distribution of X_2 given $X_1 = x_1$

$$P_{X_2|X_1}(x_2|x_1) = \frac{P_{X_1X_2}(x_1, x_2)}{P_{X_1}(x_1)}.$$

The conditional expectation: given $X_2 = x_2$

$$\mathbb{E}_{X_1 \sim P_{X_1|X_2}} [f(X_1)] = \int f(x_1) P_{X_1|X_2}(x_1|x_2) dx_1.$$

- **Bayes' theorem:** Relates both conditionals

$$P_{X_1|X_2}(x_1|x_2) = \frac{P_{X_2|X_1}(x_2|x_1)P_{X_1}(x_1)}{P_{X_2}(x_2)}.$$

Let $\mathcal{S} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots\}$ be a set of images of shape $N \times N$. For example:



- ▶ Each image $\mathbf{x}^{(i)}$ can be seen as a realization of an N^2 dimensional random variable \mathbf{X} .
- ▶ Each pixel is a random variable X_i , for $1 \leq i \leq N^2$.
- ▶ We can calculate pixelwise mean and standard deviation of pixel intensity.

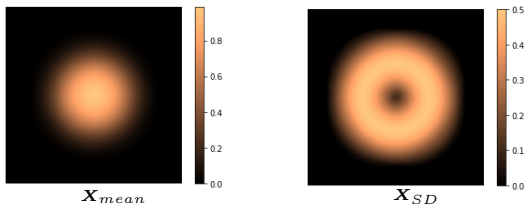
- ▶ Calculating $\mathbb{E}_{\mathbf{X} \sim P_{\mathbf{X}}} [f(\mathbf{X})] = \int f(\mathbf{x}) P_{\mathbf{X}}(\mathbf{x}) d\mathbf{x}$ may be non-trivial.
- ▶ Approximate the integral using a suitable **quadrature**.
- ▶ Monte Carlo is one such strategy:
 - ▶ Choose i.i.d. samples $\{\mathbf{x}^{(i)}\}_{i=1}^N$
 - ▶ Calculate the estimator

$$E_N[f] = \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}^{(i)}) \approx \mathbb{E}_{\mathbf{X} \sim P_{\mathbf{X}}} [f(\mathbf{X})]$$

- ▶ Error = $\mathcal{O}(\frac{1}{\sqrt{N}})$

- ▶ Useful when all we have is data

For example, with $N = 10,000$ random circle image



Consider the phenomena of heat flow through some material.

Modelled using the **heat equation**:

$$\frac{\partial u(\mathbf{s}, t)}{\partial t} - \kappa \Delta u(\mathbf{s}, t) = 0$$

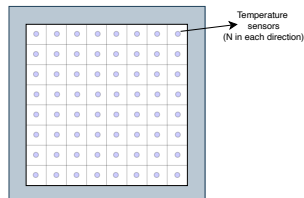
where

$$u(\mathbf{s}, 0) = u_0(\mathbf{s})$$

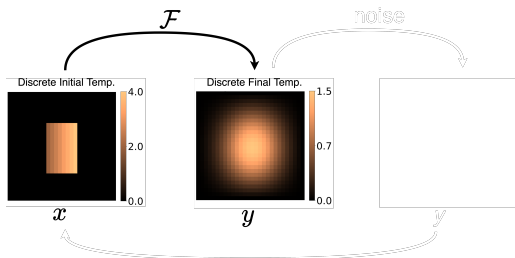
$u(\mathbf{s}, t) \rightarrow$ temperature at location \mathbf{s} at time t

$u_0(\mathbf{s}) \rightarrow$ initial temperature at location \mathbf{s}

$\kappa \rightarrow$ thermal conductivity of material



Forward problem \mathcal{F} : Given $u_0(\mathbf{s})$ at the sensor nodes determine $u(\mathbf{s}, T)$



Consider the phenomena of heat flow through some material.

Modelled using the **heat equation**:

$$\frac{\partial u(\mathbf{s}, t)}{\partial t} - \kappa \Delta u(\mathbf{s}, t) = 0$$

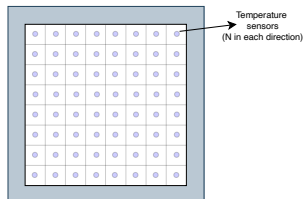
where

$$u(\mathbf{s}, 0) = u_0(\mathbf{s})$$

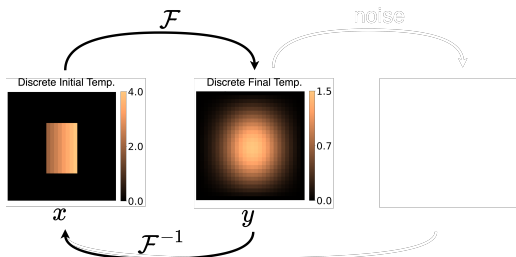
$u(\mathbf{s}, t) \rightarrow$ temperature at location \mathbf{s} at time t

$u_0(\mathbf{s}) \rightarrow$ initial temperature at location \mathbf{s}

$\kappa \rightarrow$ thermal conductivity of material



Inverse problem \mathcal{F}^{-1} : Given $u(\mathbf{s}, T)$ at the sensor nodes infer $u_0(\mathbf{s})$



Consider the phenomena of heat flow through some material.

Modelled using the **heat equation**:

$$\frac{\partial u(\mathbf{s}, t)}{\partial t} - \kappa \Delta u(\mathbf{s}, t) = 0$$

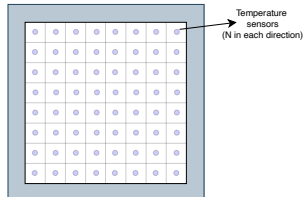
where

$$u(\mathbf{s}, 0) = u_0(\mathbf{s})$$

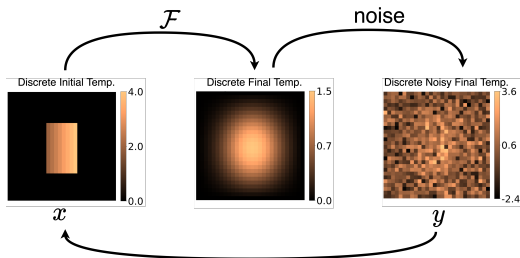
$u(\mathbf{s}, t) \rightarrow$ temperature at location \mathbf{s} at time t

$u_0(\mathbf{s}) \rightarrow$ initial temperature at location \mathbf{s}

$\kappa \rightarrow$ thermal conductivity of material



Inverse problem \mathcal{F}^{-1} : Given **noisy** $u(\mathbf{s}, T)$ infer $u_0(\mathbf{s})$

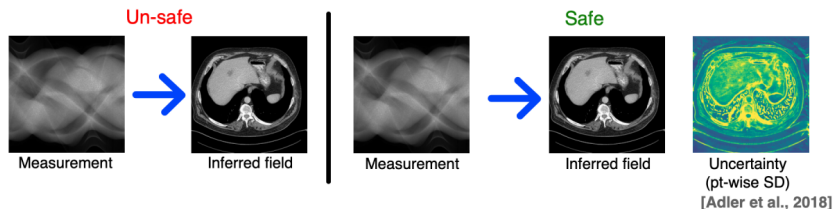


Challenges with inverse problems:

- ▶ Inverse map is **not well posed**.
- ▶ **Noisy measurements**.
- ▶ Need to encode **prior knowledge** about the inferred field x (e.g. initial temperature).

Uncertainty in inferred field critical for applications with high-stake decisions.

Example: Medical imaging to detect liver lesions



Assume x and y are modelled by random variables X and Y .

AIM: Given a measurement $Y = y$ approximate the conditional distribution

$$P_{X|Y}(x|y) = \frac{P_{Y|X}(y|x)P_X(x)}{P_Y(y)}$$

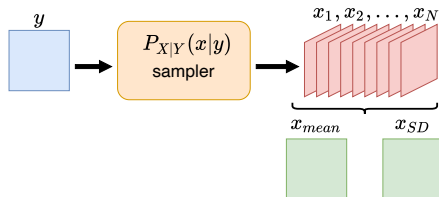
$P_X(x)$: prior distribution

$P_Y(y)$: evidence

$P_{Y|X}(y|x)$: likelihood of observing y given x

$P_{X|Y}(x|y)$: posterior distribution given y

We want to generate samples:



Bayesian formulation: challenges

- ▶ Posterior sampling techniques, such as Markov Chain Monte Carlo, are prohibitively expensive when **dimension of \mathbf{X} is large**.
- ▶ Characterization of **priors for complex data**

For example, \mathbf{x} data might look like:



Representing this data using simple distributions is **hard**!

Resolve both issues using deep learning

A neural network is a parametrized mapping

$$NN_{\theta} : \Omega_x \rightarrow \Omega_y$$

typically formed by alternating composition

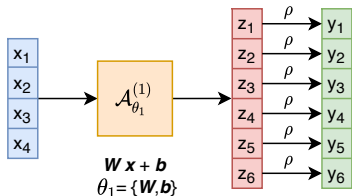
$$NN_{\theta} := \rho \circ \mathcal{A}_{\theta_{L+1}}^{(L+1)} \circ \rho \circ \mathcal{A}_{\theta_L}^{(L)} \circ \rho \circ \mathcal{A}_{\theta_{L-1}}^{(L-1)} \circ \dots \circ \rho \circ \mathcal{A}_{\theta_1}^{(1)}$$

where

$\theta = \{\theta_k\}_{k=1}^L \rightarrow$ trainable weights and biases of the network

$\mathcal{A}_{\theta_k}^{(k)} \rightarrow$ parametrized affine transformation

$\rho \rightarrow$ non-linear activation function



A neural network is a parametrized mapping

$$NN_{\theta} : \Omega_x \rightarrow \Omega_y$$

typically formed by alternating composition

$$NN_{\theta} := \rho \circ \mathcal{A}_{\theta_{L+1}}^{(L+1)} \circ \rho \circ \mathcal{A}_{\theta_L}^{(L)} \circ \rho \circ \mathcal{A}_{\theta_{L-1}}^{(L-1)} \circ \dots \circ \rho \circ \mathcal{A}_{\theta_1}^{(1)}$$

where

$\theta = \{\theta_k\}_{k=1}^L \rightarrow$ trainable weights and biases of the network

$\mathcal{A}_{\theta_k}^{(k)} \rightarrow$ parametrized affine transformation

$\rho \rightarrow$ non-linear activation function

Usage:

- ▶ Let x and y are related in some manner, say $y = f(x)$.
- ▶ We are only given $\mathcal{S} = \{(x_i, y_i)\}_{i=1}^N$.
- ▶ NN_{θ} can be used to learn f .

Consider a suitable measure

$$\mu : \Omega_x \times \Omega_y \times \Omega_x \times \Omega_y \rightarrow \mathbb{R}$$

s.t. $\mu(\mathbf{x}, \mathbf{y}, \mathbf{x}, NN_{\theta}(\mathbf{x}))$ is the error/discrepancy between $(\mathbf{x}, \mathbf{y}) \in \mathcal{S}$ and $(\mathbf{x}, NN_{\theta}(\mathbf{x}))$.

Define the loss/objective function

$$\Pi(\theta) = \frac{1}{N} \sum_{i=1}^N \mu(\mathbf{x}_i, \mathbf{y}_i, \mathbf{x}_i, NN_{\theta}(\mathbf{x}_i))$$

Solve the optimization problem

$$\theta^* = \arg \min_{\theta} \Pi(\theta)$$

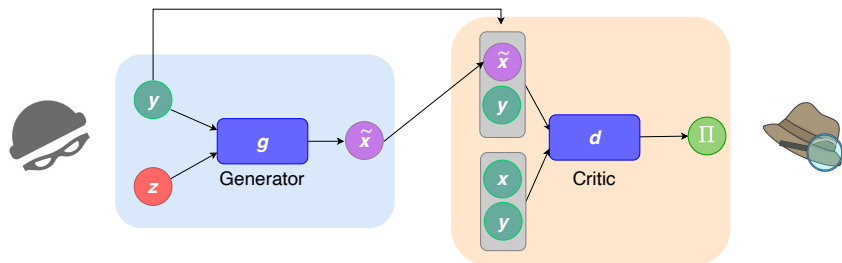
Then $NN_{\theta^*} \approx \mathbf{f}$

Also need to tune network [hyper-parameters](#):

- Width
- Depth (L)
- Activation function ρ
- Optimizer
- Loss function
- Dataset

Conditional generative adversarial network (cGAN)

- ▶ Learning distributions conditioned on another field.
- ▶ Comprises two neural networks, g and d .
- ▶ Flipping role of x and y for inverse problems.



Generator network:

- ▶ $g : \Omega_z \times \Omega_y \rightarrow \Omega_x$.
- ▶ Latent variable $\mathbf{Z} \sim P_Z$, e.g. $N(0, \mathbf{I})$. Also $N_z \ll N_x$.
- ▶ (x, y) sampled from true P_{XY}

Critic network:

- ▶ $d : \Omega_x \times \Omega_y \rightarrow \mathbb{R}$.
- ▶ d tries to detect fake samples.
- ▶ $d(x, y)$ large for real x , small otherwise.

- ▶ Given $Y = y$ and $Z \sim P_Z$ we get a random variable

$$X^g = g(Z, y), \quad X^g \sim P_{X|Y}^g.$$

- ▶ Objective function

$$\Pi(g, d) = \mathbb{E}_{\substack{(X, Y) \sim P_{XY} \\ Z \sim P_Z}} [d(X, Y) - d(g(Z, Y), Y)]$$

- ▶ g and d determined (with constraint $\|d\|_{\text{Lip}} \leq 1$) through

$$(g^*, d^*) = \arg \min_g \arg \max_d \Pi(g, d)$$

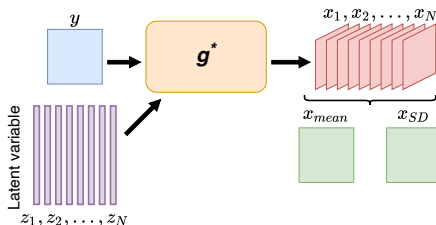
- ▶ Adler et al. (2018) proved that the minmax problem is equivalent to

$$g^* = \arg \min_g \mathbb{E}_{Y \sim P_Y} [W_1(P_{X|Y}, P_{X|Y}^g)]$$

where W_1 is the Wasserstein-1 distance.

Steps:

- ▶ Acquire samples $\{x_1, \dots, x_N\}$, where x_i are sampled from P_X .
- ▶ For each x_i acquire the measurement y_i (using \mathcal{F}). y_i will be a sample from $P_{Y|X}$
- ▶ Construct the paired set $\mathcal{S} = \{(x_1, y_1), \dots, (x_N, y_N)\}$. (x_i, y_i) can be seen as sampled from true joint P_{XY} .
- ▶ Train a cGAN on \mathcal{S} .
- ▶ For a new test measurement y , generate samples using g^* .
- ▶ Evaluate statistics using Monte Carlo



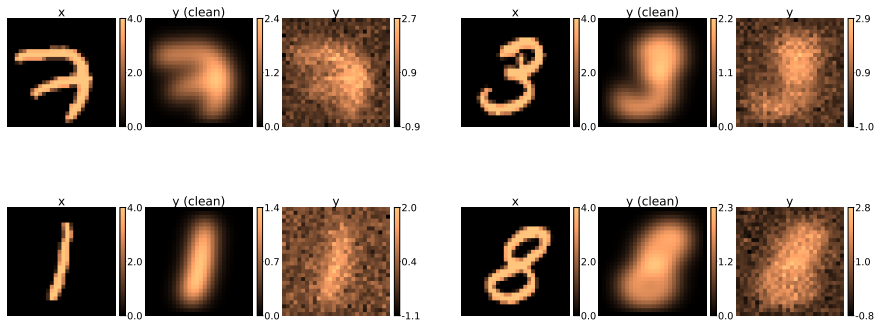
- ▶ We want to solve the inverse heat conduction problem.
- ▶ Assume the space has $N_x = 28 \times 28 = 784$ sensors.
- ▶ \mathbf{X} : Initial temperature at the sensors.
- ▶ \mathbf{Y} : Noisy final temperature at the sensors.
- ▶ We assume $\kappa = 0.2$
- ▶ (\mathbf{x}, \mathbf{y}) pairs obtained by numerically solving the forward problem \mathcal{F} and artificially adding noise.
- ▶ We need to assume some prior distribution on \mathbf{X} .

The efficacy and generalizability of conditional GANs for posterior inference in physics-based inverse problems (D. Ray, D. Patel, H. Ramaswamy, A. A. Oberai); preprint 2022.

Inverse heat equation: MNIST prior on \mathbf{X}

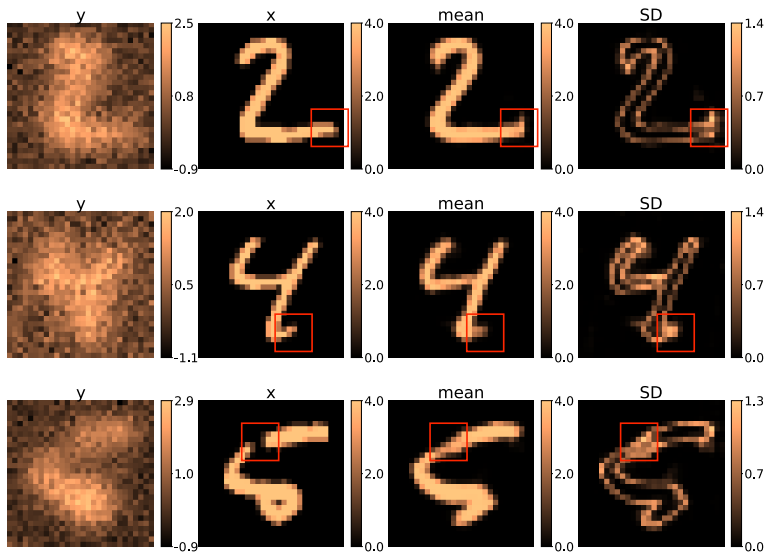
Assuming \mathbf{X} to be given by MNIST handwritten digits: u_0 takes the value 4.0 on the digit and 0.0 everywhere else.

Training samples:



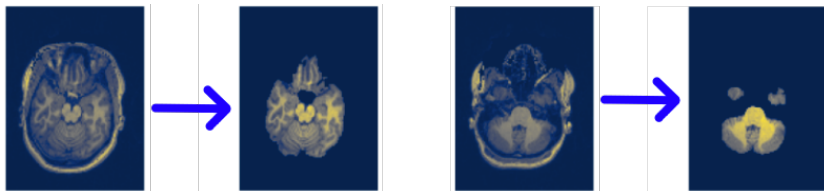
- ▶ Trained a network with dimension $N_z = 100$ for latent variable \mathbf{Z} . Note that $N_x = 784$.
- ▶ We don't have clean y in real problems!

Testing trained generator g^* . Monte Carlo statistics with 800 z samples.



MRI brain extraction problem

- ▶ MRI images of the head used in several downstream task:
 - ▶ Quantifying grey and white matter
 - ▶ Monitoring neurological diseases, e.g., Alzheimer's
 - ▶ Estimating brain atrophy
- ▶ MRI images need to undergo "skull stripping": eliminating everything other than the brain.
- ▶ Manual process can be tedious and subjective.
- ▶ Existing algorithms do not quantify uncertainty.



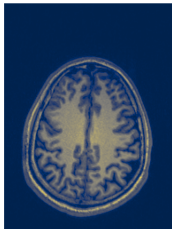
Setup

- ▶ X : Image of the brain (what we infer)
- ▶ Y : Image of the head: original MRI
- ▶ Each image has shape $256 \times 192 \implies 49,152$ pixels
- ▶ Publicly available NFBS paired dataset used for training (no explicit \mathcal{F})
- ▶ cGAN trained with latent dimension $N_z = 256$.

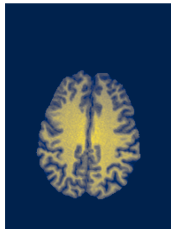
Probabilistic Brain Extraction in MR Images via Conditional Generative Adversarial Networks (A. Moazami, D. Ray, D. Pelletier, A. A. Oberai); preprint 2022.

Testing on a new sample

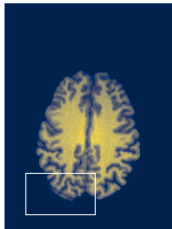
Input image



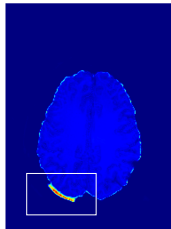
Target image



Output (mean)



Standard Deviation



- ▶ Conditional GANs can be used to approximate and sample from conditional distributions.
- ▶ Required paired data to train – supervised learning model.
- ▶ What do we gain?
 - ▶ Ability to represent and encode complex data.
 - ▶ Dimension reduction since $N_z \ll N_x$.

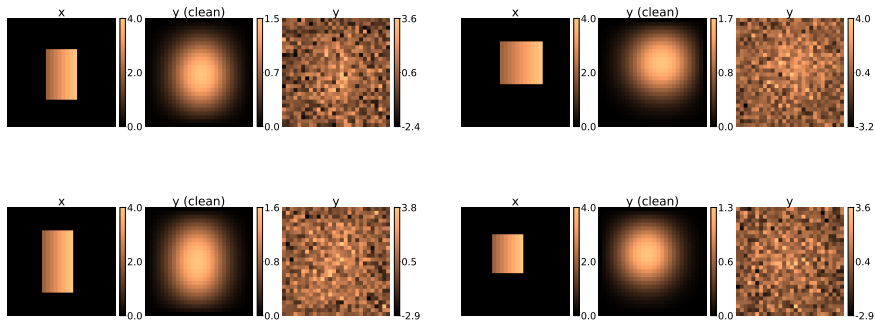
	N_x	N_z	Dim. compression
Inverse heat equation	784	100	7.84
Brain extraction	49,152	256	192

- ▶ Once trained, sampling from cGAN is quick and easy.
- ▶ Algorithm tested for many other physics-based and medical applications.

Questions?

Assuming \mathbf{X} to given by a rectangular inclusion and $N_x = N_y = 28 \times 28 = 784$

Training samples:



We never actually have clean \mathbf{y} !

Testing trained cGAN (statistics with 800 z samples)

