# Solving Inverse Problems in Imaging
## by Joint Posterior Maximization with Autoencoding Prior

Andrés Almansa

Université de Paris

Cnrs

<u>Joint work with</u>:      Mario González, Pauline Tan
<u>Thanks to</u>:    Pablo Musé, Mauricio Delbracio, José Lezama
Preprint and code available here
http://up5.fr/jpmap

March 9th 2021 - Lecture 8 - M2 MVA course
https://delires.wp.imt.fr

Andrés Almansa    Solving Inverse Problems in Imaging

Introduction | Inverse problems in Imaging
Proposed Method | Implicitly decoupled methods
Experiments | Explicitly decoupled methods

## Inverse Problems in Imaging

Estimate clean image $\boldsymbol{x} \in \mathbb{R}^d$
from noisy, degraded measurements $\boldsymbol{y} \in \mathbb{R}^m$.



Measurements $\boldsymbol{y}$      Inverse problem      Ideal image $\boldsymbol{x}$

Known degradation model (usually log-concave):

$$p_{Y|X}\left(\boldsymbol{y} \mid \boldsymbol{x}\right) \propto e^{-F(\boldsymbol{x},\boldsymbol{y})} \quad \text{where} \quad F(\boldsymbol{x},\boldsymbol{y}) = \frac{1}{2\sigma^2}\|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{y}\|^2. \quad (1)$$

Andrés Almansa    Solving Inverse Problems in Imaging

Introduction
Proposed Method
Experiments

Inverse problems in Imaging
Implicitly decoupled methods
Explicitly decoupled methods

# Inverse Problems in Imaging

Estimate clean image $\boldsymbol{x} \in \mathbb{R}^d$
from noisy, degraded measurements $\boldsymbol{y} \in \mathbb{R}^m$.
Known degradation model (usually log-concave):

$$p_{Y|X}\left(\boldsymbol{y} \mid \boldsymbol{x}\right) \propto e^{-F(\boldsymbol{x},\boldsymbol{y})} \quad \text{where} \quad F(\boldsymbol{x},\boldsymbol{y}) = \frac{1}{2\sigma^2}\|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{y}\|^2. \quad (1)$$

## Variational/Bayesian Approach

Use image prior $p_X\left(\boldsymbol{x}\right) \propto e^{-\lambda R(\boldsymbol{x})}$ to compute estimator

$$\hat{\boldsymbol{x}}_{\mathrm{MAP}} = \arg\max_{\boldsymbol{x}} p_{X|Y}\left(\boldsymbol{x} \mid \boldsymbol{y}\right) = \arg\min_{\boldsymbol{x}} \left\{F(\boldsymbol{x},\boldsymbol{y}) + \lambda R(\boldsymbol{x})\right\} \quad (2)$$

$$\hat{\boldsymbol{x}}_{\mathrm{MMSE}} = \arg\min_{\boldsymbol{x}} \mathbb{E}\left[\left. \|X - \boldsymbol{x}\|^2 \right| Y = \boldsymbol{y}\right] \quad (3)$$

Introduction | Inverse problems in Imaging
Proposed Method | Implicitly decoupled methods
Experiments | Explicitly decoupled methods

# Inverse Problems in Imaging

$$p_{Y|X}\left(\boldsymbol{y}\,|\,\boldsymbol{x}\right) \propto e^{-F(\boldsymbol{x},\boldsymbol{y})} \quad \text{where} \quad F(\boldsymbol{x},\boldsymbol{y}) = \frac{1}{2\sigma^2}\|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{y}\|^2. \quad (1)$$

### Variational/Bayesian Approach

Use image prior $p_X\left(\boldsymbol{x}\right) \propto e^{-\lambda R(\boldsymbol{x})}$ to compute estimator

$$\hat{\boldsymbol{x}}_{\text{MAP}} = \arg\max_{\boldsymbol{x}} p_{X|Y}\left(\boldsymbol{x}\,|\,\boldsymbol{y}\right) = \arg\min_{\boldsymbol{x}} \left\{F(\boldsymbol{x},\boldsymbol{y}) + \lambda R(\boldsymbol{x})\right\} \quad (2)$$

$$\hat{\boldsymbol{x}}_{\text{MMSE}} = \arg\min_{\boldsymbol{x}} \mathbb{E}\left[\left.\|X - \boldsymbol{x}\|^2\,\right|\,Y = \boldsymbol{y}\right] \quad (3)$$

Common explicit priors

- **Total Variation** (CHAMBOLLE, 2004; LOUCHET AND MOISAN, 2013; PEREYRA, 2016; RUDIN ET AL., 1992)
- **Gaussian Mixtures** (TEODORO ET AL., 2018; YU ET AL., 2011; ZORAN AND WEISS, 2011)

Introduction
Proposed Method
Experiments

Inverse problems in Imaging
Implicitly decoupled methods
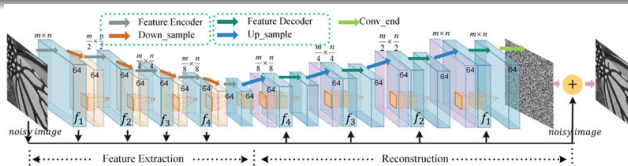Explicitly decoupled methods

Neural Networks for inverse problems:

*Two paradigms*



- Agnostic approach : find a sufficient number of image pairs $(x^i, y^i)$ and train a neural network $f_\theta$ to invert $A$ by minimizing the empirical risk $\sum_i \|f_\theta(y^i) - x^i\|_2^2$

  ✓ no need to model $A$, $\mathbf{n}$ nor prior for $x$
  ✗ needs retraining if $A$ or $\mathbf{n}$ change

Introduction
Proposed Method
Experiments

Inverse problems in Imaging
Implicitly decoupled methods
Explicitly decoupled methods

Neural Networks for inverse problems:

*Two paradigms*



- Agnostic approach : find a sufficient number of image pairs $(\boldsymbol{x}^i, \boldsymbol{y}^i)$ and train a neural network $f_\theta$ to invert $A$ by minimizing the empirical risk $\sum_i \|f_\theta(\boldsymbol{y}^i) - \boldsymbol{x}^i\|_2^2$

    ✓ no need to model $A$, $\mathbf{n}$ nor prior for $\boldsymbol{x}$

    ✗ needs retraining if $A$ or $\mathbf{n}$ change

- Decoupled (plug & play) approach : Model separately

    1. conditional density $p_{Y|X}(\boldsymbol{y} \mid \boldsymbol{x})$

        (using physical model, calibration)

    2. prior model $p_X(\boldsymbol{x})$          (through NN learning)

    3. Use Bayes theorem to estimate $\boldsymbol{x}$ via MAP or MMSE

Introduction
Proposed Method
Experiments

Inverse problems in Imaging
Implicitly decoupled methods
Explicitly decoupled methods

Neural Networks for inverse problems:

*Two paradigms*

- Agnostic approach : find a sufficient number of image pairs $(\mathbf{x}^i, \mathbf{y}^i)$ and train a neural network $f_\theta$ to invert $A$ by minimizing the empirical risk $\sum_i \|f_\theta(\mathbf{y}^i) - \mathbf{x}^i\|_2^2$

  - ✓ no need to model $A$, $\mathbf{n}$ nor prior for $\mathbf{x}$
  - ✗ needs retraining if $A$ or $\mathbf{n}$ change

- Decoupled (plug & play) approach : Model separately

  1. conditional density $p_{Y|X}(\mathbf{y}\,|\,\mathbf{x})$
     (using physical model, calibration)
  2. prior model $p_X(\mathbf{x})$ (through NN learning)
  3. Use Bayes theorem to estimate $\mathbf{x}$ via MAP or MMSE

     - ✓ uses all available modeling information
     - ✓ train once, use for many inverse problems
     - ⚠ difficult to learn $p_X(\mathbf{x})$ directly
     - ⚠ Non-convex optimization

Introduction
Proposed Method
Experiments

Inverse problems in Imaging
Implicitly decoupled methods
Explicitly decoupled methods

Neural Networks for inverse problems:

*Implicitly decoupled approach*

Solve the optimization problem

$$\hat{\boldsymbol{x}}_{\mathrm{MAP}} = \arg\max_{\boldsymbol{x}} p_{X|Y}(\boldsymbol{x} \mid \boldsymbol{y}) = \arg\min_{\boldsymbol{x}} \{F(\boldsymbol{x}, \boldsymbol{y}) + \lambda R(\boldsymbol{x})\}$$

via ADMM splitting (Ryu et al., 2019)

1. $\boldsymbol{v}_{k+1} = \arg\min_{\boldsymbol{v}} R(\boldsymbol{v}) + \frac{1}{2\delta^2}\|\boldsymbol{v} - (\boldsymbol{x}_k - \boldsymbol{u}_k)\|^2$

2. $\boldsymbol{x}_{k+1} = \arg\min_{\boldsymbol{x}} F(\boldsymbol{x}, \boldsymbol{y}) + \frac{\lambda}{2\delta^2}\|\boldsymbol{x} - (\boldsymbol{v}_{k+1} - \boldsymbol{u}_k)\|^2$

3. $\boldsymbol{u}_{k+1} = \boldsymbol{u}_k + \boldsymbol{v}_{k+1} - \boldsymbol{x}_{k+1}$

$R$ is unknown but we can use a train a neural network to approximate the $\delta$-denoising problem in step 1:

$$D_\delta(\tilde{\boldsymbol{x}}) = \arg\min_{\boldsymbol{v}} R(\boldsymbol{v}) + \frac{1}{2\delta^2}\|\boldsymbol{v} - \tilde{\boldsymbol{x}}\|^2$$

Introduction
Proposed Method
Experiments

Inverse problems in Imaging
Implicitly decoupled methods
Explicitly decoupled methods

Neural Networks for inverse problems:

*Implicitly decoupled approach*

Solve the optimization problem via ADMM splitting

$$\hat{\boldsymbol{x}}_{\mathrm{MAP}} = \arg \max_{\boldsymbol{x}} p_{X|Y}(\boldsymbol{x} \mid \boldsymbol{y}) = \arg \min_{\boldsymbol{x}} \{F(\boldsymbol{x}, \boldsymbol{y}) + \lambda R(\boldsymbol{x})\}$$

$R$ is unknown but a NN approximates its proximal operator:

$$D_\delta(\tilde{\boldsymbol{x}}) = \arg \min_{\boldsymbol{v}} R(\boldsymbol{v}) + \frac{1}{2\delta^2} \|\boldsymbol{v} - \tilde{\boldsymbol{x}}\|^2$$

Challenges

- NN training produces an MMSE rather than a MAP estimator for $D_\delta$
- Convergence guarantees

Introduction
Proposed Method
Experiments

Inverse problems in Imaging
Implicitly decoupled methods
Explicitly decoupled methods

# Neural Networks for inverse problems:

## *Implicitly decoupled approach*

Solve the optimization problem via ADMM splitting (RYU ET AL., 2019)

$$\hat{\boldsymbol{x}}_{\text{MAP}} = \arg\max_{\boldsymbol{x}} p_{X|Y}(\boldsymbol{x} \mid \boldsymbol{y}) = \arg\min_{\boldsymbol{x}} \{F(\boldsymbol{x}, \boldsymbol{y}) + \lambda R(\boldsymbol{x})\}$$

### Assumption (A)

1. ✓ $D_\delta - I$ is L-Lipschitz with $L \in (0, 1)$
2. ✗ $F(\cdot, \boldsymbol{y})$ is $\mu$-strongly convex
3. ✗ $\lambda < \frac{\sigma^2 \mu (1 + L - 2L^2)}{L} \underset{L \to 1^-}{\to} 0$

### Theorem (RYU ET AL. (2019))

*Under assumption A, the Plug & Play ADMM algorithm converges to a critical point.*

Introduction
Proposed Method
Experiments

Inverse problems in Imaging
Implicitly decoupled methods
Explicitly decoupled methods

# **Explicitly** decoupled approach (MAP-$x$):
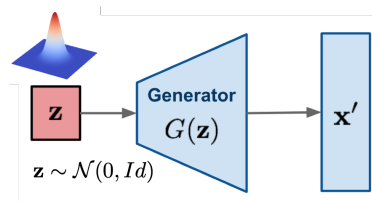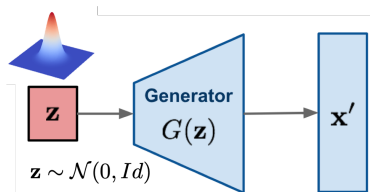
*How to use neural networks to learn the prior $p_X(x)$ ?*

Generative Adversarial Networks (GANs) (ARJOVSKY AND BOTTOU, 2017; GOODFELLOW ET AL., 2014)

Learn a generator function G that maps

$$z \sim \mathcal{N}(0, Id)$$

to

$$x = \mathsf{G}(z) \sim p_X$$

Introduction
Proposed Method
Experiments

Inverse problems in Imaging
Implicitly decoupled methods
**Explicitly decoupled methods**

## **Explicitly** decoupled approach (MAP-$\boldsymbol{x}$):
### *How to use neural networks to learn the prior $p_X(\boldsymbol{x})$ ?*

Generative Adversarial Networks (GANs) (ARJOVSKY AND BOTTOU, 2017; GOODFELLOW ET AL., 2014)

Learn a generator function G that maps

$$\boldsymbol{z} \sim \mathcal{N}(0, Id)$$

to

$$\boldsymbol{x} = \mathsf{G}(\boldsymbol{z}) \sim p_X$$



MAP-$\boldsymbol{x}$ Following PAPAMAKARIOS ET AL. (2019, SECTION 5), the push-forward measure $p_X = \mathsf{G}\sharp p_Z$ can be developed as

$$p_X(\boldsymbol{x}) = \frac{p_Z\left(\mathsf{G}^{-1}(\boldsymbol{x})\right)}{\sqrt{\det S(\mathsf{G}^{-1}(\boldsymbol{x}))}} \delta_{\mathcal{M}}(\boldsymbol{x})$$

where

$$S = \left(\frac{\partial \mathsf{G}}{\partial \boldsymbol{z}}\right)^T \left(\frac{\partial \mathsf{G}}{\partial \boldsymbol{z}}\right)$$

$$\mathcal{M} = \{\boldsymbol{x} \,:\, \exists \boldsymbol{z}, \, \boldsymbol{x} = \mathsf{G}(\boldsymbol{z})\}$$

Introduction
Proposed Method
Experiments

Inverse problems in Imaging
Implicitly decoupled methods
Explicitly decoupled methods

**Explicitly** decoupled approach (MAP-*x*):

*How to use neural networks to learn the prior $p_X(x)$ ?*

Generative Adversarial Networks (GANs) (ARJOVSKY AND BOTTOU, 2017; GOODFELLOW ET AL., 2014)  Learn a generator function G that maps

$$z \sim \mathcal{N}(0, Id) \quad \text{to} \quad x = G(z) \sim p_X$$

MAP-*x* Following PAPAMAKARIOS ET AL. (2019, SECTION 5), the push-forward measure $p_X = G\sharp p_Z$ can be developed as

$$p_X(x) = \frac{p_Z(G^{-1}(x))}{\sqrt{\det S(G^{-1}(x))}} \delta_{\mathcal{M}}(x)$$

where

$$S = \left(\frac{\partial G}{\partial z}\right)^T \left(\frac{\partial G}{\partial z}\right)$$

$$\mathcal{M} = \{x : \exists z, x = G(z)\}$$

*x*-optimization required to obtain $\hat{x}_{\mathrm{MAP}}$ becomes intractable due to:

- computation of $S$ and $\det S$,
- inversion of G, and
- hard constraint $x \in \mathcal{M}$

8/36

Andrés Almansa    Solving Inverse Problems in Imaging

Introduction
Proposed Method
Experiments

Inverse problems in Imaging
Implicitly decoupled methods
Explicitly decoupled methods

## **Explicitly** decoupled approach ($\mathrm{MAP}$-$\boldsymbol{z}$):

Instead of solving the $\boldsymbol{x}$-optimisation problem:

$$\hat{\boldsymbol{x}}_{\mathrm{MAP}} = \arg\max_{\boldsymbol{x}} p_{Y|X}\left(\boldsymbol{y} \mid \boldsymbol{x}\right) p_X\left(\boldsymbol{x}\right) = \arg\min_{\boldsymbol{x}} \left\{F(\boldsymbol{x}, \boldsymbol{y}) + R(\boldsymbol{x})\right\}$$

$\mathrm{BORA}$ $\mathrm{ET}$ $\mathrm{AL.}$ $(2017)$ propose to optimize over $\boldsymbol{z}$

$$\hat{\boldsymbol{z}} = \arg\max_{\boldsymbol{z}} \left\{p_{Y|X}\left(\boldsymbol{y} \mid \mathsf{G}(\boldsymbol{z})\right) p_Z\left(\boldsymbol{z}\right)\right\}$$

$$= \arg\min_{\boldsymbol{z}} \left\{F(\mathsf{G}(\boldsymbol{z}), \boldsymbol{y}) + \frac{1}{2}\|\boldsymbol{z}\|^2\right\}$$

$$\hat{\boldsymbol{x}}_{\boldsymbol{z}-\mathrm{MAP}} = \mathsf{G}(\hat{\boldsymbol{z}})$$

$\hat{\boldsymbol{x}}_{\boldsymbol{z}-\mathrm{MAP}}$ ($\neq \hat{\boldsymbol{x}}_{\mathrm{MAP}}$) but it maximizes the latent posterior:

$$\hat{\boldsymbol{x}}_{\boldsymbol{z}-\mathrm{MAP}} = \mathsf{G}\left(\arg\max_{\boldsymbol{z}} \left\{p_{Z|Y}\left(\boldsymbol{z} \mid \boldsymbol{y}\right)\right\}\right)$$

Introduction | Inverse problems in Imaging
Proposed Method | Implicitly decoupled methods
Experiments | Explicitly decoupled methods

**Explicitly** decoupled approach (MAP-$z$):

$\hat{x}_{z-\mathrm{MAP}}$ ($\neq \hat{x}_{\mathrm{MAP}}$) maximizes the latent posterior:

$$\hat{x}_{z-\mathrm{MAP}} = \mathsf{G}\left(\arg\max_{z}\left\{p_{Z|Y}\left(z \mid y\right)\right\}\right)$$

$$= \mathsf{G}\left(\arg\min_{z}\left\{F(\mathsf{G}(z), y) + \frac{1}{2}\|z\|^2\right\}\right)$$

### Challenges

- Nonconvex optimization using gradient descent
- may get stuck in spurious local minima

Common solution: Splitting + continuation scheme

Introduction
Proposed Method
Experiments

Inverse problems in Imaging
Implicitly decoupled methods
Explicitly decoupled methods

## MAP-$z$ splitting and continuation scheme.

$$\hat{\boldsymbol{x}}_\beta = \arg\min_{\boldsymbol{x}} \min_{\boldsymbol{z}} \underbrace{\left\{ F(\boldsymbol{x}, \boldsymbol{y}) + \frac{\beta}{2}\|\boldsymbol{x} - \mathsf{G}(\boldsymbol{z})\|^2 + \frac{1}{2}\|\boldsymbol{z}\|^2 \right\}}_{J_{1,\beta}(\boldsymbol{x}, \boldsymbol{z})}$$

$$\hat{\boldsymbol{x}}_{\mathrm{MAP}-\boldsymbol{z}} = \lim_{\beta \to \infty} \hat{\boldsymbol{x}}_\beta.$$

---

**Algorithm 1.1** MAP-$z$ splitting

**Require:** Measurements $\boldsymbol{y}$, Initial condition $\boldsymbol{x}_0$
**Ensure:** $\hat{\boldsymbol{x}} = \mathsf{G}\left(\arg\max_{\boldsymbol{z}} p_{Z|Y}(\boldsymbol{z}\,|\,\boldsymbol{y})\right)$
 1: **for** $k := 0$ **to** $k_{\max}$ **do**
 2:   $\beta := \beta_k$
 3:   **for** $n := 0$ **to** maxiter **do**
 4:     $\boldsymbol{z}_{n+1} := \arg\min_{\boldsymbol{z}} J_{1,\beta}(\boldsymbol{x}_n, \boldsymbol{z})$      // Nonconvex
 5:     $\boldsymbol{x}_{n+1} := \arg\min_{\boldsymbol{x}} J_{1,\beta}(\boldsymbol{x}, \boldsymbol{z}_{n+1})$      // Quadratic
 6:   **end for**
 7:   $\boldsymbol{x}_0 := \boldsymbol{x}_{n+1}$
 8: **end for**
 9: **return** $\boldsymbol{x}_{n+1}$

---

Non-convex step 4: Use a local quadratic approximation (VAE encoder) ...

Introduction
**Proposed Method**
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
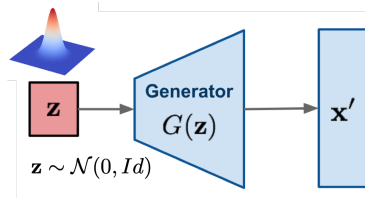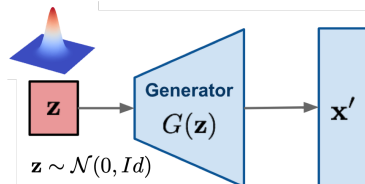Denoising Criterion and Continuation Scheme

# VAEs and Joint Posterior

**Generative Adversarial Networks (GANs)** (GOODFELLOW ET AL., 2014)

Learn a generator function G that maps

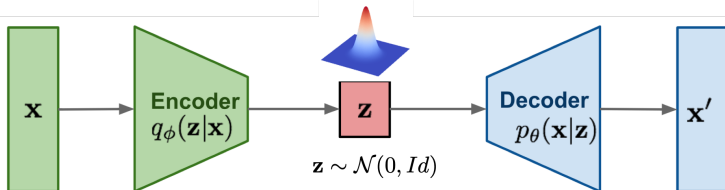$$z \sim \mathcal{N}(0, Id)$$

to

$$x = \mathsf{G}(z) \sim p_X$$

Introduction
**Proposed Method**
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
Denoising Criterion and Continuation Scheme

# VAEs and Joint Posterior

**Generative Adversarial Networks (GANs)** (GOODFELLOW ET AL., 2014)

Learn a generator function G that maps

$$z \sim \mathcal{N}(0, Id)$$

to

$$x = \mathsf{G}(z) \sim p_X$$



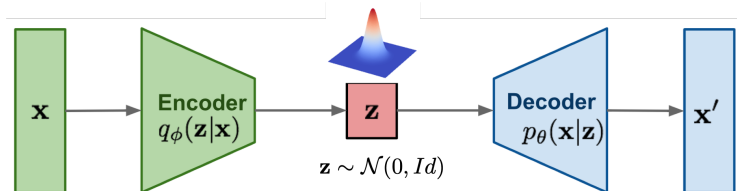**Variational AutoEncoders (VAEs)** (KINGMA AND WELLING, 2013)



Generative model: $\quad p_{X|Z}\left(x \mid z\right) = p_\theta(x|z) = \mathcal{N}(x; \boldsymbol{\mu}_\theta(z), \gamma Id)$

Approximate inverse: $\quad p_{Z|X}\left(z \mid x\right) \approx q_\phi(z|x) = \mathcal{N}(z; \boldsymbol{\mu}_\phi(x), \boldsymbol{\Sigma}_\phi(x))$

Introduction
**Proposed Method**
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
Denoising Criterion and Continuation Scheme

# VAEs and Joint Posterior

**Variational AutoEncoders (VAEs)** (KINGMA AND WELLING, 2013)



$\mathbf{z} \sim \mathcal{N}(0, Id)$

Generative model: $\qquad p_{X|Z}(\boldsymbol{x} \mid \boldsymbol{z}) = p_\theta(\boldsymbol{x}|\boldsymbol{z}) = \mathcal{N}(\boldsymbol{x}; \boldsymbol{\mu}_\theta(\boldsymbol{z}), \gamma Id)$

Approximate inverse: $\qquad p_{Z|X}(\boldsymbol{z} \mid \boldsymbol{x}) \approx q_\phi(\boldsymbol{z}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{z}; \boldsymbol{\mu}_\phi(\boldsymbol{x}), \boldsymbol{\Sigma}_\phi(\boldsymbol{x}))$

Learning: Maximize the averaged *Evidence Lower BOund (ELBO)* for $\boldsymbol{x} \in \mathcal{D}$

$$\mathcal{L}_{\theta,\phi}(\boldsymbol{x}) = \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}[\log p_\theta(\boldsymbol{x}|\boldsymbol{z})] - KL(q_\phi(\boldsymbol{z}|\boldsymbol{x}) \mid\mid p_Z(\boldsymbol{z})) \leq \log p_\theta(\boldsymbol{x}).$$

Introduction
**Proposed Method**
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
Denoising Criterion and Continuation Scheme

# VAEs and Joint Posterior

**Variational AutoEncoders (VAEs)** (Kingma and Welling, 2013)

Generative model: $\quad p_{X|Z}(x \mid z) = p_\theta(x|z) = \mathcal{N}(x; \mu_\theta(z), \gamma Id)$

Joint density: $\quad p_{X,Z}(x, z) = p_\theta(x|z) \, p_Z(z)$

Approximate inverse: $\quad p_{Z|X}(z \mid x) \approx q_\phi(z|x) = \mathcal{N}(z; \mu_\phi(x), \Sigma_\phi(x))$

Approximate joint density: $\quad \tilde{p}_{X,Z}(x, z) := q_\phi(z|x) \, p_X(x) \approx p_{X,Z}(x, z)$

Introduction
**Proposed Method**
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
Denoising Criterion and Continuation Scheme

# VAEs and Joint Posterior

**Variational AutoEncoders (VAEs)** (Kingma and Welling, 2013)

Generative model: $\quad\quad\quad\quad\quad\quad p_{X|Z}(x \mid z) = p_\theta(x|z) = \mathcal{N}(x; \mu_\theta(z), \gamma Id)$

Joint density: $\quad\quad\quad\quad\quad\quad\quad\quad p_{X,Z}(x, z) = p_\theta(x|z)\, p_Z(z)$

Approximate inverse: $\quad\quad\quad\quad p_{Z|X}(z \mid x) \approx q_\phi(z|x) = \mathcal{N}(z; \mu_\phi(x), \Sigma_\phi(x))$

Approximate joint density: $\quad\quad \tilde{p}_{X,Z}(x, z) := q_\phi(z|x)\, p_X(x) \approx p_{X,Z}(x, z)$

Joint Posterior: (log-quadratic in $x$)

$$
\begin{aligned}
J_1(x, z) &:= -\log p_{X,Z|Y}(x, z \mid y) \\
&= -\log p_{Y|X,Z}(y \mid x, z)\, p_\theta(x \mid z) p_Z(z) \\
&= F(x, y) + \underbrace{\frac{1}{2\gamma}\|x - \mu_\theta(z)\|^2 + \frac{1}{2}\|z\|^2}_{H_\theta(x,z)}.
\end{aligned}
\tag{4}
$$

Introduction
**Proposed Method**
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
Denoising Criterion and Continuation Scheme

# VAEs and Joint Posterior

**Variational AutoEncoders (VAEs)** (Kingma and Welling, 2013)

Generative model: $\qquad\qquad p_{X|Z}(x \mid z) = p_\theta(x|z) = \mathcal{N}(x; \mu_\theta(z), \gamma Id)$

Joint density: $\qquad\qquad\qquad\qquad p_{X,Z}(x,z) = p_\theta(x|z)\, p_Z(z)$

Approximate inverse: $\qquad p_{Z|X}(z \mid x) \approx q_\phi(z|x) = \mathcal{N}(z; \mu_\phi(x), \Sigma_\phi(x))$

Approximate joint density: $\qquad \tilde{p}_{X,Z}(x,z) := q_\phi(z|x)\, p_X(x) \approx p_{X,Z}(x,z)$

Joint Posterior: (log-quadratic in $x$)

$$
\begin{aligned}
J_1(x,z) &:= -\log p_{X,Z|Y}(x,z \mid y) \\
&= -\log p_{Y|X,Z}(y \mid x,z)\, p_\theta(x \mid z) p_Z(z) \\
&= F(x,y) + \underbrace{\frac{1}{2\gamma}\|x - \mu_\theta(z)\|^2 + \frac{1}{2}\|z\|^2}_{H_\theta(x,z)}.
\end{aligned} \tag{4}
$$

Approximate Joint Posterior: (log-quadratic in $z$)

$$
\begin{aligned}
J_2(x,z) &:= -\log p_{Y|X,Z}(y \mid x,z)\, q_\phi(z \mid x)\, p_X(x) \\
&= F(x,y) + \underbrace{\frac{1}{2}\|\Sigma_\phi^{-1/2}(x)(z - \mu_\phi(x))\|^2 + C(x)}_{K_\phi(x,z)} - \log p_X(x).
\end{aligned} \tag{5}
$$

Introduction
Proposed Method
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
Denoising Criterion and Continuation Scheme

## Joint Posterior Maximization - Alternate Convex Search

---

**Algorithm 2.1** Joint posterior maximization - exact case

**Require:** Measurements $\boldsymbol{y}$, Autoencoder parameters $\theta$, $\phi$, Initial condition $\boldsymbol{x}_0$

**Ensure:** $\hat{\boldsymbol{x}}, \hat{\boldsymbol{z}} = \arg\max_{\boldsymbol{x},\boldsymbol{z}} p_{X,Z|Y}(\boldsymbol{x}, \boldsymbol{z} \mid \boldsymbol{y})$

1: **for** $n := 0$ **to** maxiter **do**
2:     $\boldsymbol{z}_{n+1} := \arg\min_{\boldsymbol{z}} J_2(\boldsymbol{x}_n, \boldsymbol{z}) = \boldsymbol{\mu}_\phi(\boldsymbol{x}_n)$   // Quadratic approx
3:     $\boldsymbol{x}_{n+1} := \arg\min_{\boldsymbol{x}} J_1(\boldsymbol{x}, \boldsymbol{z}_{n+1})$             // Quadratic
4: **end for**
5: **return** $\boldsymbol{x}_{n+1}, \boldsymbol{z}_{n+1}$

---

### Proposition

If the encoder approximation is exact ($J_2 = J_1$) then

- $J_1$ is biconvex, and following GORSKI ET AL. (2007):
- Algorithm 2.1 is an Alternate Convex Search
- Algorithm 2.1 converges to a critical point

Introduction
Proposed Method
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
Denoising Criterion and Continuation Scheme

## JPMAP - Accuracy of encoder approximation



Contour plots of $-\log p_{Z|X}(z \mid x)$ and $-\log q_\phi(z|x)$ for a fixed $x$ and for a random 2D subspace in the $z$ domain.

Introduction
Proposed Method
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
Denoising Criterion and Continuation Scheme

# JPMAP - Accuracy of encoder approximation



(a) Encoder approximation   (b) Decoded exact optimum   (c) Decoded approx. optimum   (d) Difference (b)-(c)

*Figure 1. Encoder approximation: (a)* Contour plots of $-\log p_\theta(x|z) + \frac{1}{2}\|z\|^2$ and $-\log q_\phi(z|x)$ for a fixed $x$ and for a random 2D subspace in the $z$ domain (the plot shows $\pm 2\Sigma_\phi^{1/2}$ around $\mu_\phi$). Observe the relatively small gap between the true posterior $p_\theta(z|x)$ and its variational approximation $q_\phi(z|x)$. This figure shows some evidence of partial $z$-convexity of $J_1$ around the minimum of $J_2$, but it does not show how far is $z^1$ from $z^2$. *(b)* Decoded exact optimum $x_1 = \mu_\theta \left( \arg\max_z p_\theta(x|z)e^{\frac{1}{2}\|z\|^2} \right)$. *(c)* Decoded approximate optimum $x_2 = \mu_\theta \left( \arg\max_z q_\phi(z|x) \right)$. *(d)* Difference betweeen (b) and (c)

Introduction
Proposed Method
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
Denoising Criterion and Continuation Scheme

## Joint Posterior Maximization - approximate case

---

**Algorithm 2.2** Joint posterior maximization - approximate case

---

**Require:** Measurements $\boldsymbol{y}$, Autoencoder parameters $\theta$, $\phi$, Initial conditions $\boldsymbol{x}_0, \boldsymbol{z}_0$

**Ensure:** $\hat{\boldsymbol{x}}, \hat{\boldsymbol{z}} = \arg\max_{\boldsymbol{x}, \boldsymbol{z}} p_{X,Z|Y}(\boldsymbol{x}, \boldsymbol{z} \mid \boldsymbol{y})$

1: **for** $n := 0$ **to** maxiter **do**
2:     $\boldsymbol{z}^1 := \arg\min_{\boldsymbol{z}} J_2(\boldsymbol{x}_n, \boldsymbol{z}) = \boldsymbol{\mu}_\phi(\boldsymbol{x}_n)$    // Quadratic approx
3:     $\boldsymbol{z}^2 := \text{GD}_{\boldsymbol{z}} J_1(\boldsymbol{x}_n, \boldsymbol{z})$, starting from $\boldsymbol{z} = \boldsymbol{z}^1$
4:     $\boldsymbol{z}^3 := \text{GD}_{\boldsymbol{z}} J_1(\boldsymbol{x}_n, \boldsymbol{z})$, starting from $\boldsymbol{z} = \boldsymbol{z}_n$
5:     **for** $i := 1$ **to** $3$ **do**
6:         $\boldsymbol{x}^i := \arg\min_{\boldsymbol{x}} J_1(\boldsymbol{x}, \boldsymbol{z}^i)$               // Quadratic
7:     **end for**
8:     $i^* := \arg\min_{i \in \{1,2,3\}} J_1(\boldsymbol{x}^i, \boldsymbol{z}^i)$
9:     $(\boldsymbol{x}_{n+1}, \boldsymbol{z}_{n+1}) := (\boldsymbol{x}^{i^*}, \boldsymbol{z}^{i^*})$
10: **end for**
11: **return** $\boldsymbol{x}_{n+1}, \boldsymbol{z}_{n+1}$

---

Introduction
Proposed Method
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
Denoising Criterion and Continuation Scheme

## Joint Posterior Maximization - approximate case

**Algorithm 2.3** Joint posterior maximization - approximate case (faster version)

**Require:** Measurements $\boldsymbol{y}$, Autoencoder parameters $\theta$, $\phi$, Initial condition $\boldsymbol{x}_0$, iterations $n_1 \le n_2 \le n_{\max}$

**Ensure:** $\hat{\boldsymbol{x}}, \hat{\boldsymbol{z}} = \arg\max_{\boldsymbol{x},\boldsymbol{z}} p_{X,Z|Y}(\boldsymbol{x}, \boldsymbol{z} \mid \boldsymbol{y})$

1: **for** $n := 0$ **to** $n_{\max}$ **do**
2:      done := FALSE
3:      **if** $n < n_1$ **then**
4:          $\boldsymbol{z}^1 := \arg\min_{\boldsymbol{z}} J_2(\boldsymbol{x}_n, \boldsymbol{z}) = \boldsymbol{\mu}_\phi(\boldsymbol{x}_n)$      // Quadratic approx
5:          $\boldsymbol{x}^1 := \arg\min_{\boldsymbol{x}} J_1(\boldsymbol{x}, \boldsymbol{z}^1)$      // Quadratic
6:          **if** $J_1(\boldsymbol{x}^1, \boldsymbol{z}^1) < J_1(\boldsymbol{x}_n, \boldsymbol{z}_n)$ **then**
7:              $i^* := 1$      // Faster alternative while $J_2$ is good enough
8:              done := TRUE
9:          **end if**
10:      **end if**
11:      **if not** done **and** $n < n_2$ **then**
12:          $\boldsymbol{z}^1 := \arg\min_{\boldsymbol{z}} J_2(\boldsymbol{x}_n, \boldsymbol{z}) = \boldsymbol{\mu}_\phi(\boldsymbol{x}_n)$      // Quadratic approx
13:          $\boldsymbol{z}^2 := \text{GD}_{\boldsymbol{z}} J_1(\boldsymbol{x}_n, \boldsymbol{z})$, starting from $\boldsymbol{z} = \boldsymbol{z}^1$
14:          $\boldsymbol{x}^2 := \arg\min_{\boldsymbol{x}} J_1(\boldsymbol{x}, \boldsymbol{z}^2)$      // Quadratic
15:          **if** $J_1(\boldsymbol{x}^2, \boldsymbol{z}^2) < J_1(\boldsymbol{x}_n, \boldsymbol{z}_n)$ **then**
16:              $i^* := 2$      // $J_2$ init is good enough
17:              done := TRUE
18:          **end if**
19:      **end if**
20:      **if not** done **then**
21:          $\boldsymbol{z}^3 := \text{GD}_{\boldsymbol{z}} J_1(\boldsymbol{x}_n, \boldsymbol{z})$, starting from $\boldsymbol{z} = \boldsymbol{z}_n$
22:          $\boldsymbol{x}^3 := \arg\min_{\boldsymbol{x}} J_1(\boldsymbol{x}, \boldsymbol{z}^3)$      // Quadratic
23:          $i^* := 3$
24:      **end if**
25:      $(\boldsymbol{x}_{n+1}, \boldsymbol{z}_{n+1}) := (\boldsymbol{x}^{i^*}, \boldsymbol{z}^{i^*})$
26: **end for**
27: **return** $\boldsymbol{x}_{n+1}, \boldsymbol{z}_{n+1}$

Introduction
Proposed Method
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
Denoising Criterion and Continuation Scheme

# JPMAP - Effectivenes of the encoder initialization

Trajectories of $\mathrm{GD}_{\boldsymbol{z}} J_1(\boldsymbol{x}_0, \boldsymbol{z})$, starting from $\boldsymbol{z} = \boldsymbol{z}_0$
Thick blue curve: $\boldsymbol{z}_0 = \arg\min_{\boldsymbol{z}} J_2(\boldsymbol{x}_0, \boldsymbol{z}) = \boldsymbol{\mu}_\phi(\boldsymbol{x}_0)$
Thin curves: random initializations $\boldsymbol{z}_0 \sim \mathcal{N}(0, Id)$

Introduction
Proposed Method
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
Denoising Criterion and Continuation Scheme

# JPMAP - Convergence

If we use ELU activations then the following assumption is verified:

## Assumption (2)

$J_1(\cdot, \boldsymbol{z})$ is convex and admits a minimizer for any $\boldsymbol{z}$. Moreover, $J_1$ is coercive and continuously differentiable.

## Proposition (Convergence of Algorithm 2.3)

Let $\{(\boldsymbol{x}_n, \boldsymbol{z}_n)\}$ be a sequence generated by Algorithm 2.3. Under Assumption 2 we have that:

1. The sequence $\{J_1(\boldsymbol{x}_n, \boldsymbol{z}_n)\}$ converges monotonically when $n \to \infty$.

2. The sequence $\{(\boldsymbol{x}_n, \boldsymbol{z}_n)\}$ has at least one accumulation point.

3. All accumulation points of $\{(\boldsymbol{x}_n, \boldsymbol{z}_n)\}$ are stationary points of $J_1$ and they all have the same function value.

Introduction
**Proposed Method**
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
**Denoising Criterion and Continuation Scheme**

# Denoising Criterion to train VAEs (IM ET AL., 2017)

**Variational AutoEncoders (VAEs)** (KINGMA AND WELLING, 2013)



$$\mathbf{z} \sim \mathcal{N}(0, Id)$$

Learning: Maximize the averaged *Evidence Lower BOund (ELBO)* for $\boldsymbol{x} \in \mathcal{D}$

$$\mathcal{L}_{\theta,\phi}(\boldsymbol{x}) = \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}[\log p_\theta(\boldsymbol{x}|\boldsymbol{z})] - KL(q_\phi(\boldsymbol{z}|\boldsymbol{x}) \parallel p_Z(\boldsymbol{z})) \leq \log p_\theta(\boldsymbol{x}).$$

Introduction
**Proposed Method**
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
**Denoising Criterion and Continuation Scheme**

# Denoising Criterion to train VAEs (IM ET AL., 2017)

**Variational AutoEncoders (VAEs)** (KINGMA AND WELLING, 2013)



Learning: Maximize the averaged *Evidence Lower BOund (ELBO)* for $\boldsymbol{x} \in \mathcal{D}$

$$\mathcal{L}_{\theta,\phi}(\boldsymbol{x}) = \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}[\log p_\theta(\boldsymbol{x}|\boldsymbol{z})] - KL(q_\phi(\boldsymbol{z}|\boldsymbol{x}) \,||\, p_Z(\boldsymbol{z})) \leq \log p_\theta(\boldsymbol{x}).$$

Problem: $\boldsymbol{\mu}_\phi(\boldsymbol{x})$ only trained for $\boldsymbol{x} \in \mathcal{D}$ or $\boldsymbol{x} \in \mathcal{M} = \boldsymbol{\mu}_\theta(\mathbb{R}^m)$.
**But:** Step 2 in the algorithm evaluates $\boldsymbol{\mu}_\phi(\boldsymbol{x}_n)$ for degraded $\boldsymbol{x}_n \notin \mathcal{M}$

Introduction
**Proposed Method**
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
**Denoising Criterion and Continuation Scheme**

# Denoising Criterion to train VAEs (IM ET AL., 2017)

**Variational AutoEncoders (VAEs)** (KINGMA AND WELLING, 2013)



Learning: Maximize the averaged *Evidence Lower BOund (ELBO)* for $x \in \mathcal{D}$

$$\mathcal{L}_{\theta,\phi}(x) = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - KL(q_\phi(z|x) \,||\, p_Z(z)) \leq \log p_\theta(x).$$

Denoising criterion: Train on $\tilde{\mathcal{D}}$ but still require $\boldsymbol{\mu}_\theta(\boldsymbol{\mu}_\phi(\tilde{x})) \approx x$.

$$\tilde{\mathcal{D}} = \{\tilde{x} = x + \sigma_{\text{DVAE}}\varepsilon \,:\, x \in \mathcal{D} \text{ and } \varepsilon \sim \mathcal{N}(0, I)\}$$

Maximize the denoising ELBO

$$\tilde{\mathcal{L}}_{\theta,\phi}(x) = \mathbb{E}_{p(\tilde{x}|x)}\left[\mathbb{E}_{q_\phi(z|\tilde{x})}[\log p_\theta(x|z)] - KL(q_\phi(z|\tilde{x}) \,||\, p_Z(z))\right]$$

Introduction
**Proposed Method**
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
**Denoising Criterion and Continuation Scheme**

# Denoising criterion does not degrade generative model



(a) Denoising          (b) Compressed Sensing          (c) Inpainting

**Figure 1.** *Evaluating the quality of the generative model as a function of $\sigma_{\mathrm{DVAE}}$. On (a) Denoising (Gaussian noise $\sigma = 150$), (b) Compressed Sensing ($\sim 10.2\%$ measurements, noise $\sigma = 10$) and (c) Inpainting (80% of missing pixels, noise $\sigma = 10$). Results of both algorithms are computed on a batch of 50 images and initialising on ground truth $\boldsymbol{x}^*$ (for CSGM we use $\boldsymbol{z}_0 = \boldsymbol{\mu}_\phi(\boldsymbol{x}^*)$).*

Introduction
**Proposed Method**
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
**Denoising Criterion and Continuation Scheme**

# Optimal value of $\sigma_{\mathrm{DVAE}}$



(a) Denoising

(b) Compressed Sensing

(c) Inpainting

**Figure 2.** *Evaluating the effectiveness of JPMAP vs CGSM as a function of $\sigma_{\mathrm{DVAE}}$ (same setup of Figure 1). Without a denoising criterion $\sigma_{\mathrm{DVAE}} = 0$ the JPMAP algorithm may provide wrong guesses $z^1$ when applying the encoder in step 2 of Algorithm 2.2. For $\sigma_{\mathrm{DVAE}} > 0$ however, the alternating minimization algorithm can benefit from the robust initialization heuristics provided by the encoder, and it consistently converges to a better local optimum than the simple gradient descent in CSGM.*

Introduction
Proposed Method
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
Denoising Criterion and Continuation Scheme

## MAP-z as the limit case for $\beta \to \infty$

Two options for MAP-$z$ estimator instead of the joint MAP-$x$-$z$

1. CSGM - gradient descent, may be stuck in local minima
2. Use Algorithm 2.3 to solve a series of joint MAP-$x$-$z$ problems with increasing values of $\beta = \frac{1}{\gamma} \to \infty$ as suggested in Algorithm 1.1.

Stopping criterion: Inequality constrained problem

$$\underset{x,z\,:\,\|G(z)-x\|^2 \leq \varepsilon}{\arg\min} F(x, y) + \frac{1}{2}\|z\|^2.$$

The corresponding Lagrangian form is

$$\max_{\beta} \min_{x,z} F(x, y) + \frac{1}{2}\|z\|^2 + \beta \left( \|G(z) - x\|^2 - \varepsilon \right)^+ \qquad (6)$$

We use the exponential multiplier method (Tseng and Bertsekas, 1993) to guide the search for the optimal value of $\beta$ (see Algorithm 2.4)

Introduction
Proposed Method
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
Denoising Criterion and Continuation Scheme

# MAP-z as the limit case for $\beta \to \infty$

**Algorithm 2.4** MAP-$z$ as the limit of joint MAP-$x$-$z$.

**Require:** Measurements $y$, Tolerance $\varepsilon$, Rate $\rho > 0$, Initial $\beta_0$, Initial $x_0$, Iterations $0 \le n_1 \le n_2 \le n_{max}$

**Ensure:** $\arg\min_{z\,:\,\|G(z)-x\|^2 \le \varepsilon} F(x, y) + \frac{1}{2}\|z\|^2$.

1: $\beta := \beta_0$
2: $x^0, z^0 :=$ Algorithm 2.3 starting from $x = x_0$ with $\beta, n_1, n_2, n_{max}$.
3: converged := FALSE
4: $k := 0$
5: **while not** converged **do**
6: $\quad x^{k+1}, z^{k+1} :=$ Algorithm 2.3 starting from $x = x^k$ with $\beta$ and $n_1 = n_2 = 0$
7: $\quad C = \|G(z^{k+1}) - x^{k+1}\|^2 - \varepsilon$
8: $\quad \beta := \beta \exp(\rho C)$
9: $\quad$ converged := $(C \le 0)$
10: $\quad k := k + 1$
11: **end while**
12: **return** $x^k, z^k$

Introduction
**Proposed Method**
Experiments

Variational AutoEncoder Priors
Joint Posterior Maximization with AutoEncoding Prior
**Denoising Criterion and Continuation Scheme**

# MAP-z as the limit case for $\beta \to \infty$



**Figure 5.** Evolution of Algorithm 2.4. *In this inpainting example, JPMAP starts with the initialization in (a). During first iterations (b) − (d) where $\beta_k$ is small, $\boldsymbol{x}_k$ and $\mathsf{G}(\boldsymbol{z}_k)$ start loosely approaching each other*

# Denoising experiments (MNIST)

Andrés Almansa    Solving Inverse Problems in Imaging

# Compressed sensing experiments (MNIST)

Andrés Almansa    Solving Inverse Problems in Imaging

# Inpainting experiments (MNIST)

Andrés Almansa    Solving Inverse Problems in Imaging

# Denoising experiment: $\sigma = 110/255$



$x^*$

$y$

CSGM

JPMAP ($\beta = 1/\gamma^2$)

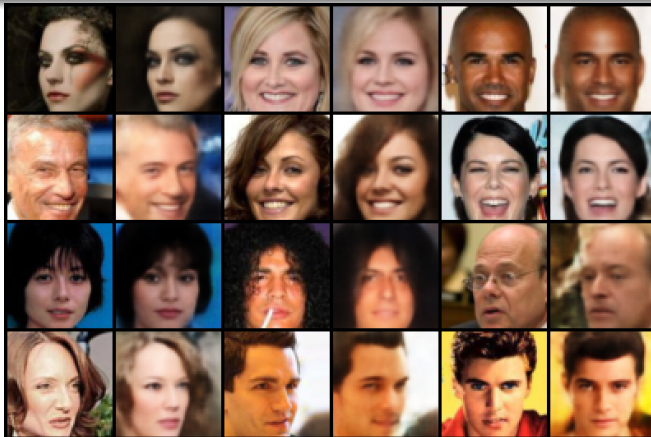JPMAP ($\alpha = 10$)

JPMAP ($\alpha = 5$)

JPMAP ($\alpha = 1$)

# Compressed sensing experiment: $m = 140$ random measurements

Andrés Almansa    Solving Inverse Problems in Imaging

# Inpainting experiment: 80% missing pixels



$x^*$

$y$

CSGM

JPMAP ($\beta = 1/\gamma^2$)

JPMAP ($\alpha = 10$)

JPMAP ($\alpha = 5$)

JPMAP ($\alpha = 1$)

# Inpainting experiment: 80% missing pixels $\sigma = 10/255$ (CelebA)



From top to bottom: original image $\boldsymbol{x}^*$, corrupted image $\tilde{\boldsymbol{x}}$, restored by CSGM, restored image $\hat{\boldsymbol{x}}$ by our framework.

# CelebA reconstructions $\mu_\theta(\mu_\phi(x))$



Reconstructions $\mu_\theta(\mu_\phi(x))$ (even columns) for some test samples $x$ (odd columns), showing the over-regularization of data manifold imposed by the trained VAE. As a consequence, $-\log p_{Z|Y}(z \mid y)$ does not have as many local minima and then a simple gradient

# Conclusion

- JPMAP avoids spurious local minima thanks to
  - Quasi bi-convex optimization
  - Encoder initialization
  - Denoising VAE
  - Splitting and continuation scheme
- JPMAP converges for all quadratic problems and regularisation parameters (unlike denoiser-based PnP approaches (RYU ET AL., 2019) that are more restrictive)
- Constraints
  - Fixed size
  - VAEs lag behind GANs

## Future work

- Use a more powerful VAE like NVAE (VAHDAT AND KAUTZ, 2020) or TwoStageVAE (DAI AND WIPF, 2019)
- Patch-based JPMAP (EPLL-like)
- Use ADMM with non-linear constraints instead of continuation scheme for MAP $- z$
- Generalize the scheme to perform posterior sampling

Andrés Almansa    Solving Inverse Problems in Imaging

Preprint and code available here
http://up5.fr/jpmap

Thank you for your attention!

Questions? Comments

📄 Arjovsky, Martin and Léon Bottou (2017). "Towards Principled Methods for Training Generative Adversarial Networks". In: *(ICLR) International Conference on Learning Representations*, pp. 1–17 (cit. on pp. 12–14).

📄 Bora, Ashish, Ajil Jalal, Eric Price, and Alexandros G Dimakis (2017). "Compressed sensing using generative models". In: *(ICML) International Conference on Machine Learning*. Vol. 2. JMLR. org, pp. 537–546. ISBN: 9781510855144. arXiv: `arXiv:1703.03208v1` (cit. on p. 15).

📄 Chambolle, A (2004). "An algorithm for total variation minimization and applications". In: *Journal of Mathematical Imaging and Vision* 20, pp. 89–97. DOI: `10.1023/B:JMIV.0000011325.36760.1e` (cit. on p. 5).

📄 Dai, Bin and David Wipf (2019). "Diagnosing and Enhancing VAE Models". In: *ICLR*, pp. 1–42. arXiv: `1903.05789` (cit. on p. 48).

📄 Goodfellow, Ian J., Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio (2014). "Generative Adversaria l Networks". In: *Advances in Neural Information Processing Systems 27*, pp. 2672–2680. ISSN: 10495258. arXiv: `1406.2661` (cit. on pp. 12–14, 18, 19).

📄 Gorski, Jochen, Frank Pfeuffer, and Kathrin Klamroth (2007). "Biconvex sets and optimization with biconvex functions: a survey and extensions". In: *Mathematical Methods of Operations Research* 66.3, pp. 373–407. ISSN: 1432-2994. DOI: `10.1007/s00186-007-0161-1` (cit. on p. 24).

Im, Daniel Jiwoong, Sungjin Ahn, Roland Memisevic, and Yoshua Bengio (2017). "Denoising criterion for variational auto-encoding framework". In: *31st AAAI Conference on Artificial Intelligence, AAAI 2017*. AAAI press, pp. 2059–2065. arXiv: 1511.06406 (cit. on pp. 31–33).

Kingma, Diederik P and Max Welling (2013). "Auto-Encoding Variational Bayes". In: *(ICLR) International Conference on Learning Representations*. MI, pp. 1–14. ISBN: 1312.6114v10. DOI: 10.1051/0004-6361/201527329. arXiv: 1312.6114 (cit. on pp. 19–23, 31–33).

Louchet, Cécile and Lionel Moisan (2013). "Posterior expectation of the total variation model: Properties and experiments". In: *SIAM Journal on Imaging Sciences* 6.4, pp. 2640–2684. ISSN: 19364954. DOI: 10.1137/120902276 (cit. on p. 5).

Papamakarios, George, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan (2019). "Normalizing Flows for Probabilistic Modeling and Inference". In: arXiv: 1912.02762 (cit. on pp. 13, 14).

Pereyra, Marcelo (2016). "Proximal Markov chain Monte Carlo algorithms". In: *Statistics and Computing* 26.4, pp. 745–760. ISSN: 0960-3174. DOI: 10.1007/s11222-015-9567-4. arXiv: 1306.0187 (cit. on p. 5).

Rudin, Leonid I., Stanley Osher, and Emad Fatemi (1992). "Nonlinear total variation based noise removal algorithms". In: *Physica D: Nonlinear Phenomena* 60.1-4, pp. 259–268. ISSN: 01672789. DOI: 10.1016/0167-2789(92)90242-F (cit. on p. 5).

Ryu, Ernest K., Jialin Liu, Sicheng Wang, Xiaohan Chen, Zhangyang Wang, and Wotao Yin (2019). "Plug-and-Play Methods Provably Converge with Properly Trained Denoisers". In: *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, pp. 5546–5557. arXiv: 1905.05406 (cit. on pp. 9, 11, 47).

Teodoro, Afonso M., José M. Bioucas-Dias, and Mário A. T. Figueiredo (2018). *Scene-Adapted Plug-and-Play Algorithm with Guaranteed Convergence: Applications to Data Fusion in Imaging*. arXiv: 1801.00605 (cit. on p. 5).

Tseng, Paul and Dimitri P. Bertsekas (1993). "On the convergence of the exponential multiplier method for convex programming". In: *Mathematical Programming* 60.1-3, pp. 1–19. ISSN: 00255610. DOI: 10.1007/BF01580598 (cit. on p. 36).

Vahdat, Arash and Jan Kautz (2020). "Nvae: A deep hierarchical variational autoencoder". In: *Advances in Neural Information Processing Systems* 33 (cit. on p. 48).

Yu, Guoshen, Guillermo Sapiro, and Stéphane Mallat (2011). "Solving inverse problems with piecewise linear estimators: From Gaussian mixture models to structured sparsity". In: *IEEE Transactions on Image Processing* 21.5, pp. 2481–2499 (cit. on p. 5).

Zoran, Daniel and Yair Weiss (2011). "From learning models of natural image patches to whole image restoration". In: *2011 International Conference on*

*Computer Vision*. IEEE, pp. 479–486. ISBN: 978-1-4577-1102-2. DOI: 10.1109/ICCV.2011.6126278 (cit. on p. 5).