

Bayesian inference with learnt generative image priors encoded by neural networks

Dr. Marcelo Pereyra

<http://www.macs.hw.ac.uk/~mp71/>

Maxwell Institute for Mathematical Sciences & Heriot-Watt University

Paris, 8 December 2022



Outline

- 1 Introduction
- 2 Bayesian imaging with generative priors supported on manifolds
- 3 Illustrative numerical experiments with a VAE prior
- 4 Scaling to high dimensions with conditional normalising flow models
- 5 Conclusion

Imaging inverse problems

- We are interested in an unknown image $x^* \in \mathbb{R}^d$.
- We measure $y \in Y$, related to x^* by some mathematical model.
- For example, in many imaging problems

$$y = Ax^* + w,$$

for some operator A that is poorly conditioned or rank deficient, and an unknown perturbation or “noise” w .

- The recovery of x^* from y is usually not well posed. Additional information is required in order to deliver meaningful solutions.

Mathematical imaging frameworks

- There are three main mathematical and computational frameworks for inference in imaging inverse problems:
 - ① Applied analysis
 - ② Bayesian statistics.
 - ③ Machine learning.
- These frameworks have complementary strengths and weaknesses.
- Our aim is a unifying framework of theory, methods, and algorithms that inherits the benefits of each approach.

The Bayesian statistical approach

- Model x^* as a realisation of a r.v. \mathbb{x} on \mathbb{R}^d . Use the distribution of \mathbb{x} to regularise the problem and promote expected properties.
- The observation y is a realisation of a r.v. $(y|\mathbb{x} = x^*)$.
- Inferences about x^* from y are derived from the joint distribution of (\mathbb{x}, y) - specified via the decomposition $p(\mathbb{x}, y) = p(y|\mathbb{x})p(\mathbb{x})$.
- This determines the posterior distribution, with density

$$p(x|y) = \frac{p(y|x)p(x)}{\int_{\mathbb{R}^d} p(y|\tilde{x})p(\tilde{x})d\tilde{x}},$$

which models our beliefs about \mathbb{x} after observing $y = y$.

- A simple algorithm to compute probabilities and expectations w.r.t. $p(x|y)$ is the Unadjusted Langevin Algorithm (ULA), given by

$$X_{k+1} = X_k + \delta_k \nabla \log p(y|X_k) + \delta_k \nabla \log p(X_k) + \sqrt{2\delta_k} Z_{k+1},$$

where $Z_{k+1} \sim \mathcal{N}(0, I_d)$ and $(\delta_k)_{k \in \mathbb{N}}$ is a sequence of step-sizes.

- The samples generated by ULA can be used to compute Monte Carlo estimates of \hat{x}_{MMSE} and perform advanced inferences.
- ULA is remarkably well understood and provably convergent under easily verifiable conditions on $p(x|y)$.

(The SDE underpinning ULA)

- Important side note: ULA arises from discrete-time approximations of the Langevin diffusion process

$$\mathbf{X}: \quad d\mathbf{X}_t = \frac{1}{2} \nabla \log p(\mathbf{X}_t|y) dt + dW_t, \quad 0 \leq t \leq T, \quad \mathbf{X}(0) = x_0.$$

where W is the Brownian motion on \mathbb{R}^d .

- When $x \mapsto p(x|y) \in \mathcal{C}^1$ with $x \mapsto \nabla \log p(x|y)$ Lipschitz continuous, X_t converges exponentially fast to $p(x|y)$ as $t \rightarrow \infty$.
- ULA stems from a basic Euler approximations of \mathbf{X} .
- We recommend to use an accelerated approximation of \mathbf{X} for significantly faster convergence (see 10.1137/19M1283719).

Outline

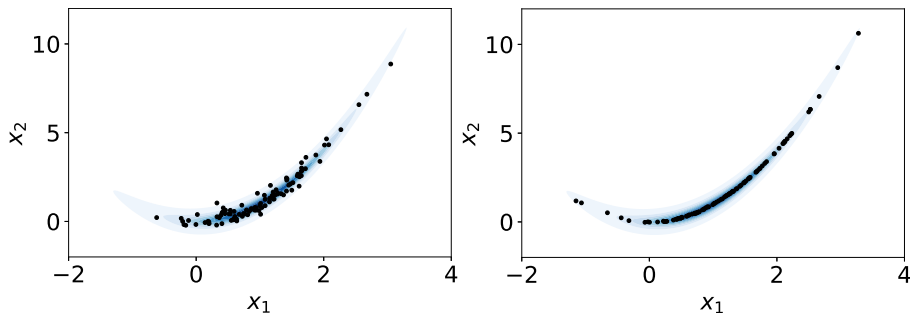
- 1 Introduction
- 2 Bayesian imaging with generative priors supported on manifolds
- 3 Illustrative numerical experiments with a VAE prior
- 4 Scaling to high dimensions with conditional normalising flow models
- 5 Conclusion

Generative image priors encoded by neural networks

Here, we focus on Bayesian inference based on deep generative priors for problems with abundant training data available to describe \mathbb{X} :

- 1 Let $\{x'_i\}_{i=1}^M$ be a training dataset that represents our prior knowledge about \mathbb{X} .
- 2 We adopt a manifold hypothesis and suppose that \mathbb{X} takes values close to an unknown p -dimensional submanifold of \mathbb{R}^d .
- 3 To estimate the manifold, we introduce a latent r.v. \mathbb{Z} on \mathbb{R}^p , with $p \ll d$, and a mapping $\nu_\theta : \mathbb{R}^p \mapsto \mathbb{R}^d$, such that the push-forward measure of $\mathbb{Z} \sim \mathcal{N}(0, I_p)$ under ν_θ is close to $\{x'_i\}_{i=1}^M$ (in dist.).
- 4 We implement ν_θ as a neural network. Can learn ν_θ from $\{x'_i\}_{i=1}^M$ by using, e.g., a VAE, a GAN, or a normalising flow approach.

Illustrative example - Rosenbrock distribution



Left: training data from the two-dimensional Rosenbrock distribution. Right: push-forward of $\mathbb{Z} \sim \mathcal{N}(0, I_p)$ under ν_θ as implemented by a VAE, with $p = 1$.

Posterior distributions for generative priors

- With \mathbb{z} and ν_θ , we have the likelihood $p(y|z) = p(y|x = \nu_\theta(z))$.
- We use Bayes' theorem to derive the posterior for $\mathbb{z}|\mathbb{y} = y$

$$p(z|y) = \frac{p(y|x = \nu_\theta(z))p(z)}{\int_{\mathbb{R}^p} p(y|\tilde{z})p(\tilde{z})d\tilde{z}},$$

- Pushing $(\mathbb{z}|\mathbb{y} = y)$ under $\nu_\theta(z)$ leads to the posterior for $(\mathbb{x}|\mathbb{y} = y)$, which supported on a manifold and does not have a density.

Key questions

Some fundamental questions:

- ① Under what conditions on the generative model are the resulting Bayesian models well-posed and amenable to efficient computation? Do the key quantities of interest inherit this well-posed nature?
- ② Are these Bayesian methods and algorithms delivering solutions that are meaningful from a non-subjective point of view?
- ③ Can we guarantee the convergence of ULA under easily verifiable conditions, with non-asymptotic accuracy bounds?

In this short talk, we will focus on the first two questions and demonstrate the approach with some numerical experiments.

For technical details please see:

- ① M. Holden, M. Pereyra, K. Zygalakis, “Bayesian Imaging with Data-Driven Priors Encoded by Neural Networks”, SIAM Journal on Imaging Sciences, 15 (2), 2022.
<https://doi.org/10.1137/21M1406313>.
- ② S. Melidonis, M. Holden, M. Pereyra, K. Zygalakis, “Empirical Bayesian imaging with conditional generative priors encoded by neural networks”, in preparation.

The oracle Bayesian model

We analyse Bayesian models with data-driven priors in an *M-complete* modelling framework:

- There exists a true - albeit unknown - marginal distribution for \mathbf{x} and posterior distribution for $(\mathbf{x}|\mathbf{y} = y)$.
- Basing inferences on these oracle models is theoretically optimal.
- We henceforth denote this optimal prior distribution by μ . When μ admits a density w.r.t. the Leb. measure on \mathbb{R}^d , we denote it by p^* .
- In that case, the posterior for $\mathbf{x}|\mathbf{y}$ has density

$$p^*(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x})p^*(\mathbf{x})}{\int_{\mathbb{R}^d} p(\mathbf{y}|\tilde{\mathbf{x}})p^*(\tilde{\mathbf{x}})d\tilde{\mathbf{x}}}.$$

The oracle Bayesian model

We analyse Bayesian models with data-driven priors in an *M-complete* modelling framework:

- Conceptual construction: μ naturally depends on the application.
- We regard the training data $\{x'_i\}_{i=1}^M$ as a sample from μ .
- When we learn ν_θ and approximate μ by assuming that $\mathbb{x} = \nu_\theta(\mathbb{z})$ for $\mathbb{z} \sim \mathcal{N}(0, \mathbf{I}_p)$, pushing $(\mathbb{z}|\mathbb{y} = y)$ under ν_θ leads to the posterior for $(\mathbb{x}|\mathbb{y} = y)$ that approximates the oracle $p^*(x|y)$.
- Holden et al. (2022a) establishes that $(\mathbb{z}|\mathbb{y} = y)$ and $(\mathbb{x}|\mathbb{y} = y)$ are well-posed in the sense of Hadamard and have finite moments.

Outline

- 1 Introduction
- 2 Bayesian imaging with generative priors supported on manifolds
- 3 Illustrative numerical experiments with a VAE prior
- 4 Scaling to high dimensions with conditional normalising flow models
- 5 Conclusion

Illustrative experiments

- We first illustrate the proposed approach with the MNIST dataset.
- We perform the following advanced inferences:
 - ① Identify the latent dimension p .
 - ② Perform MMSE inference in challenging image denoising, inpainting, and deblurring experiments.
 - ③ Adopt a likelihood-ratio test to detect out-of-sample observations that should not be analysed with the Bayesian model.
 - ④ Assess the frequentist accuracy of the Bayesian probabilities reported by the model.
- We report comparisons with MAP estimation under the same model, and with PnP-ADMM by using a DnCNN denoiser.

Identification of manifold dimension p

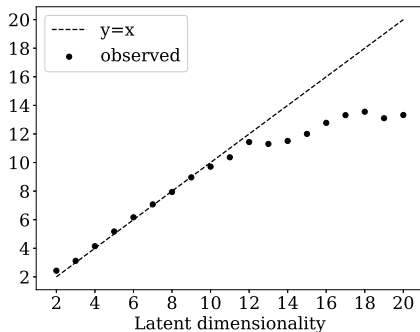
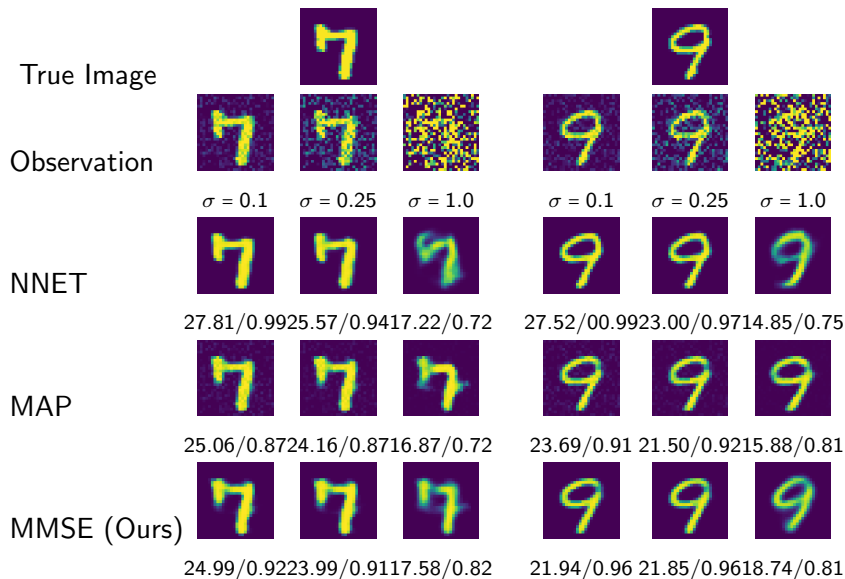
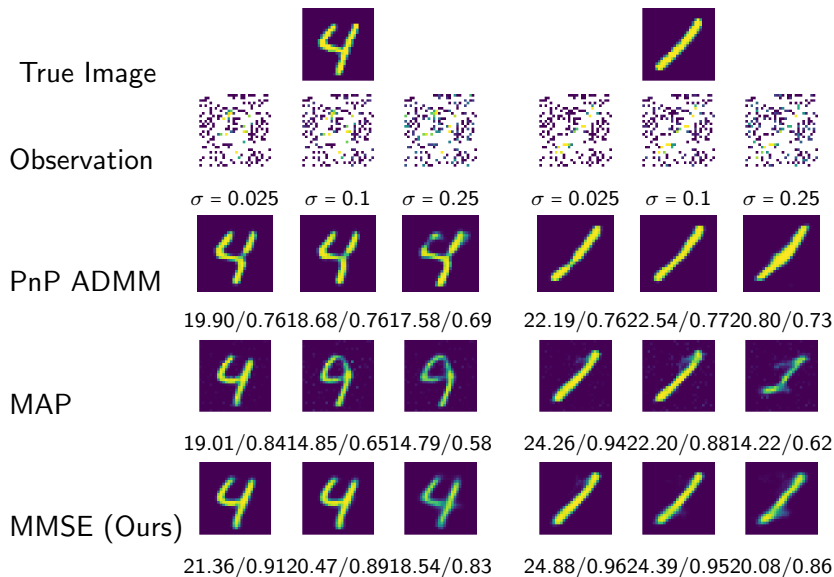


Figure: Trace of sample covariance of $\nu_{\theta}(x_i)$ across all test images. The amount of information encoded by the prior stabilises for $p \approx 12$, additional dimensions do not significantly increase the amount of prior information encoded.

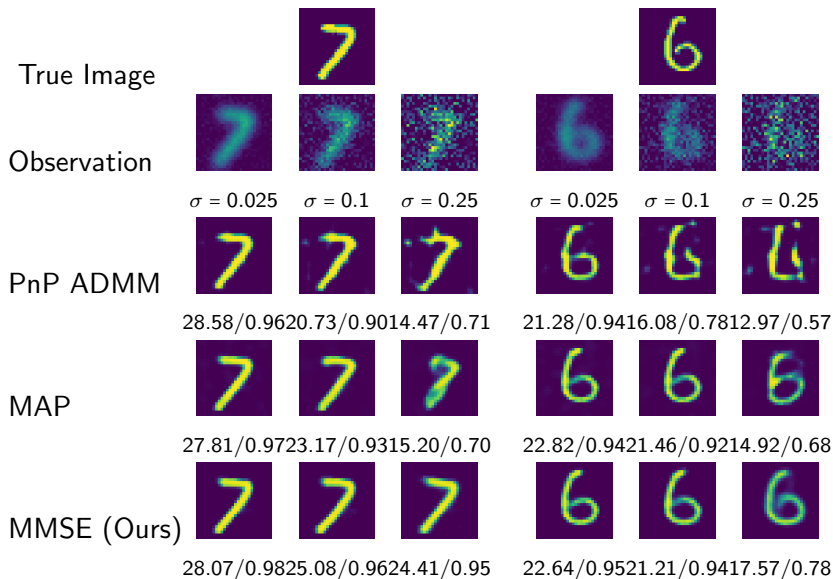
Denoising



Inpainting



Deconvolution



Likelihood ratio test for out-of-distribution detection

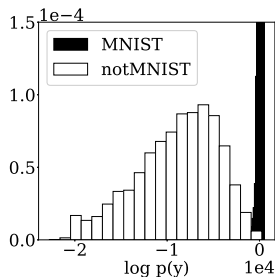


Figure: Denoising

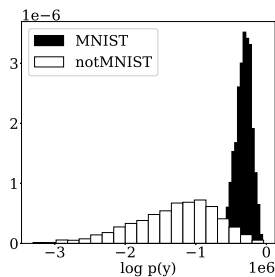


Figure: Inpainting

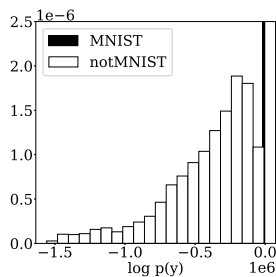


Figure: Deblurring

Figure: Histograms of marginal likelihoods for image denoising, inpainting and deblurring experiments. Out-of-sample detection powers for notMNIST of 99.6%, 88.5% and 99.8% respectively.

Coverage test: frequentist accur. of Bayesian probabilities

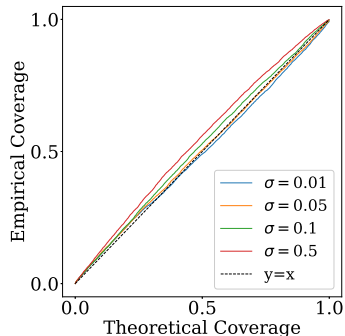


Figure: Denoising

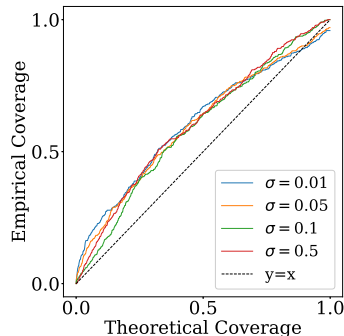
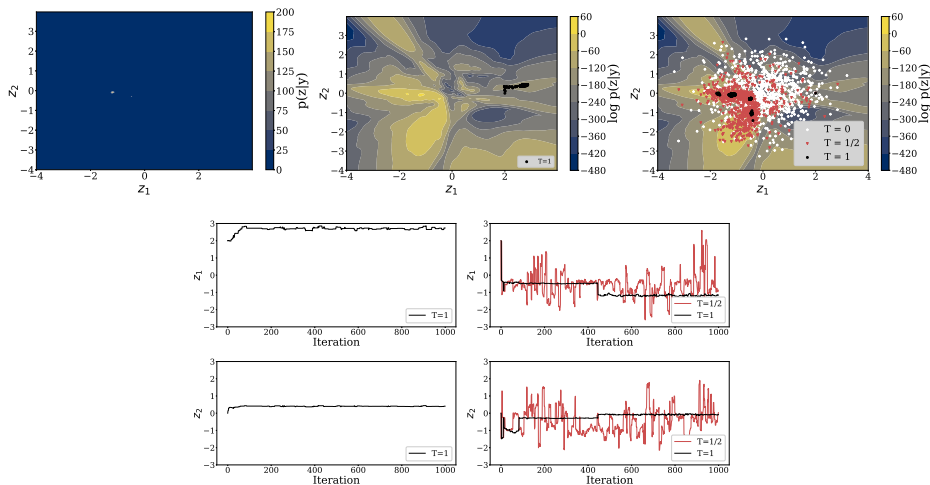


Figure: Inpainting

Multimodality? Use temperatures..

Toy example in dimension $p = 2$: three interacting Markov chains targeting $p_T(z|y) \propto p(y|z)^T p(z)$ for $T = \{0, 0.5, 1\}$.



Outline

- 1 Introduction
- 2 Bayesian imaging with generative priors supported on manifolds
- 3 Illustrative numerical experiments with a VAE prior
- 4 Scaling to high dimensions with conditional normalising flow models
- 5 Conclusion

Conditional generative priors

- Despite their success in computer vision, scaling generative models to large inference problems reliably is difficult because of mode collapse, spurious modes, or other sources of bias.
- To reduce the difficulty of the machine learning problem, we consider a **conditional generative model** $\mathbb{x} = \nu_{\theta}^u(\mathbb{z})$, $\mathbb{z} \sim \mathcal{N}(0, \mathbf{I}_p)$, that models the distribution of \mathbb{x} given some additional r.v. $\mathbb{u} = u$.
- For example, we let \mathbb{u} denote a low resolution version of \mathbb{x} , and implement ν_{θ}^u by using a normalising flow for image super-resolution.
- This leads to the model

$$p(z|y, u) = \frac{p(y|z, u)p(z)}{p(y|u)},$$

with $p(y|z, u) = p(y|x = \nu_{\theta}^u(z))$ and $p(y|u) = \int_{\mathbb{R}^p} p(y|\tilde{z}, u)p(\tilde{z})d\tilde{z}$.

Empirical Bayesian imaging with conditional generative priors

- We accurately estimate u^* from y by maximum marginal likelihood estimation:

$$\hat{u} = \underset{\mu}{\operatorname{argmax}} p_{\theta}(y|u).$$

- Adopting an empirical Bayesian strategy, we perform inference on $(\mathbb{x}|\mathbb{y} = y, \mathbb{u} = \hat{u})$ by using

$$p(z|y, \hat{u}) = \frac{p(y|z, \hat{u})p(z)}{p(y|\hat{u})},$$

and pushing to $(\mathbb{z}|\mathbb{y} = y, \mathbb{u} = \hat{u})$ via ν_{θ}^u .

- This can be performed computationally efficiently by using ULA within a Stochastic Approximation Proximal Gradient scheme that simultaneously computes \hat{u} and draws samples from $(\mathbb{z}|\mathbb{y} = y, \mathbb{u} = \hat{u})$. See <https://doi.org/10.1007/s11222-020-09986-y> and <https://doi.org/10.1137/20M1339829> for details.

Illustrative example - Image deblurring

We compare the MMSE estimator obtained with the proposed method and the MAP-style estimator obtained from PnP-ADMM, by using a DnCNN end-to-end denoising prior, on an image deblurring problem.

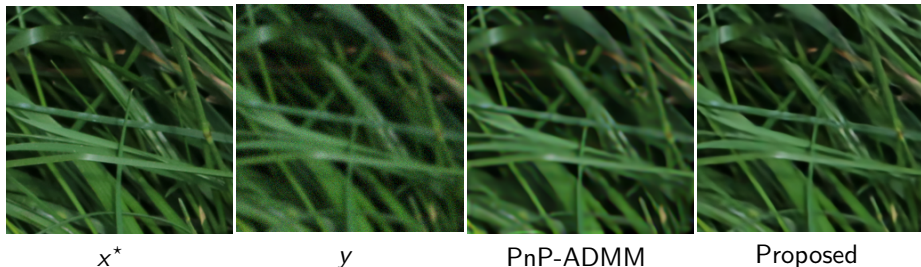
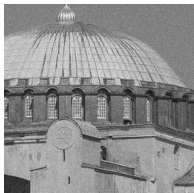


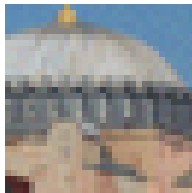
Figure: Image deblurring experiment: 9×9 uniform blur, $\sigma = 10$.

Illustrative example - Image pan-sharpening

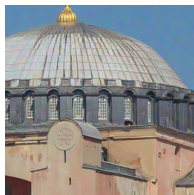
We seek to recover x^* from two noisy linear observations y_1 and y_2 , one with spectral fine details and the other with spatial fine detail.



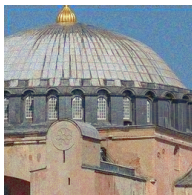
y_1



y_2



Proposed (31.5dB)



PnP ADMM (28.5dB)



x^*

Outline

- 1 Introduction
- 2 Bayesian imaging with generative priors supported on manifolds
- 3 Illustrative numerical experiments with a VAE prior
- 4 Scaling to high dimensions with conditional normalising flow models
- 5 Conclusion

Conclusion

- We have studied methodology for Bayesian inference with generative priors encoded by neural networks, learnt from training data.
- We rooted our study in the Bayesian *M-complete* paradigm that views generative Bayesian models as approximations of an oracle model.
- A key challenge to scale the approach to large problems is that generative models struggle to learn high-dimensional distributions.
- We have addressed that difficulty by adopting an empirical Bayesian approach and considering a conditional generative prior.

Thank you!