



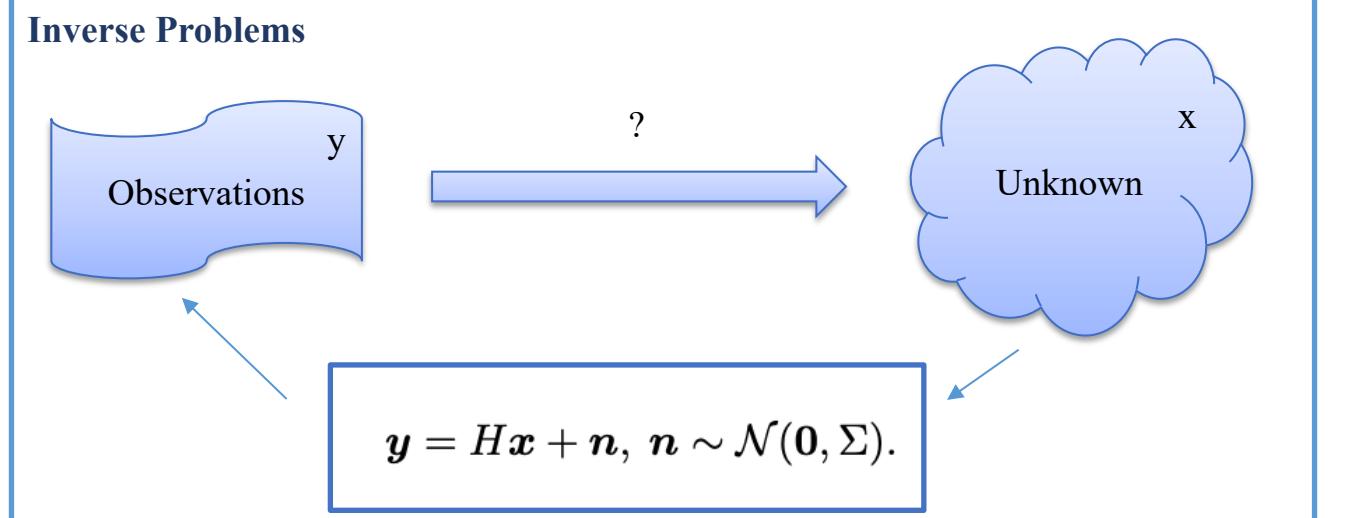
AGEM: Solving Linear Inverse Problems via Deep Priors and Sampling

Bichuan Guo¹, Yuxing Han², Jiangtao Wen¹¹Tsinghua University, ²South China Agricultural University

Abstract

In this paper we propose to use a denoising autoencoder (DAE) prior to simultaneously solve a linear inverse problem and estimate its noise parameter. Existing DAE-based methods estimate the noise parameter empirically or treat it as a tunable hyper-parameter. We instead propose *autoencoder guided EM*, a probabilistically sound framework that performs Bayesian inference with intractable deep priors. We show that efficient posterior sampling from the DAE can be achieved via Metropolis-Hastings, which allows the Monte Carlo EM algorithm to be used. We demonstrate competitive results for signal denoising, image deblurring and image devignetting. Our method is an example of combining the representation power of deep learning with uncertainty quantification from Bayesian statistics.

RELATED WORK



Discriminative learning: requires task-specific training

- Task 1 – training 1 – model 1
 - Task 2 – training 2 – model 2
- Learned knowledge is non-trivial to transfer.

Model-based learning: require analytic models

- Behavioral model: noise, blur kernel, etc.
 - Inductive bias: TV-L1, CSF, etc.
- Method fits into the Bayesian paradigm:

$$\log \Pr(\mathbf{x} | \mathbf{y}, \Sigma) = \log \Pr(\mathbf{y} | \mathbf{x}, \Sigma) + \log \Pr(\mathbf{x}) + \text{const.}$$

Drawback: requires handcrafting and often involves manual hyper-param tuning.

	Pros	Cons	Examples
analytic	tractability interpretability	bias prone	total variation
data-driven	scalability less handcrafting	intractable	databases, neural nets

A special data-driven prior with good tractability: **denoising autoencoder**

$$\mathbf{r}^*(\mathbf{x}) = \mathbf{x} + \sigma_{\text{tr}}^2 \nabla_{\mathbf{x}} \log \Pr(\mathbf{x}) + o(\sigma_{\text{tr}}^2), \text{ as } \sigma_{\text{tr}} \rightarrow 0,$$

OBJECTIVES AND METHODS

Let the data speak for itself

- Choose model-based learning over discriminative learning
 - Free ourselves from tedious re-training
- Use data-driven prior over analytic prior
 - Free ourselves from handcrafting priors
- Automatic model tuning
 - Free ourselves from tedious hyper-param searching

Determining hyper-parameters

- Perform MLE on observations.
- Clean data "x" are latent variables.
- The canonical way to perform MLE with latent variables is to use the EM (expectation maximization) algorithm.

EM is intractable

$$\begin{aligned} Q(\Sigma, \Sigma^{(\tau)}) &= \mathbb{E}_{\mathbf{x} \sim \Pr(\mathbf{x} | \mathbf{y}, \Sigma^{(\tau)})} \log \Pr(\mathbf{y}, \mathbf{x} | \Sigma) \\ &= \mathbb{E}_{\mathbf{x} \sim \Pr(\mathbf{x} | \mathbf{y}, \Sigma^{(\tau)})} \log \Pr(\mathbf{y} | \mathbf{x}, \Sigma) + \log \Pr(\mathbf{x}), \end{aligned}$$

No free lunch: we cannot expect to enjoy the advantages of data-driven priors, while demanding the tractability of analytic priors!

$$\Pr(\mathbf{x} | \mathbf{y}, \Sigma^{(\tau)}) = Z^{-1} \Pr(\mathbf{y} | \mathbf{x}, \Sigma^{(\tau)}) \Pr(\mathbf{x}),$$

Sampling methods

- Monte Carlo EM algorithm: no tractable posterior? A group of samples will also do. How to sample from the posterior?
- The first term is easy (analytic)
 - The second term is still intractable. Our only lead is its score function:

$$\mathbf{r}^*(\mathbf{x}) = \mathbf{x} + \sigma_{\text{tr}}^2 \nabla_{\mathbf{x}} \log \Pr(\mathbf{x}) + o(\sigma_{\text{tr}}^2), \text{ as } \sigma_{\text{tr}} \rightarrow 0,$$

We adopt Metropolis-Hastings algorithms to sample from the posterior:

$$\alpha = \frac{\Pr(\mathbf{x}^* | \mathbf{y}, \Sigma^{(\tau)}) q(\mathbf{x}^{(i)} | \mathbf{x}^*)}{\Pr(\mathbf{x}^{(i)} | \mathbf{y}, \Sigma^{(\tau)}) q(\mathbf{x}^* | \mathbf{x}^{(i)})},$$

$$\begin{aligned} \log \alpha &= \log \Pr(\mathbf{x}^* | \mathbf{y}, \Sigma^{(\tau)}) - \log \Pr(\mathbf{x}^{(i)} | \mathbf{y}, \Sigma^{(\tau)}) \\ &= \log \Pr(\mathbf{y} | \mathbf{x}^*, \Sigma^{(\tau)}) - \log \Pr(\mathbf{y} | \mathbf{x}^{(i)}, \Sigma^{(\tau)}) + \log \Pr(\mathbf{x}^*) - \log \Pr(\mathbf{x}^{(i)}) \\ &= \left(H \frac{\mathbf{x}^{(i)} + \mathbf{x}^*}{2} - \mathbf{y} \right)^T \Sigma^{(\tau)-1} H (\mathbf{x}^{(i)} - \mathbf{x}^*) + \log \Pr(\mathbf{x}^*) - \log \Pr(\mathbf{x}^{(i)}), \end{aligned}$$

The score function is sufficient:

$$\begin{aligned} \log \Pr(\mathbf{x}^*) - \log \Pr(\mathbf{x}^{(i)}) &\approx \nabla_{\mathbf{x}} \log \Pr(\mathbf{x})|_{\mathbf{x}^{(i)}} \cdot (\mathbf{x}^* - \mathbf{x}^{(i)}) \\ &\approx \sigma_{\text{tr}}^{-2} (r(\mathbf{x}^{(i)}) - \mathbf{x}^{(i)})^\top (\mathbf{x}^* - \mathbf{x}^{(i)}), \end{aligned}$$

To accelerate Markov chain mixing, we use Langevin-adjusted MH:

$$q_{\text{MALA}}(\mathbf{x} | \mathbf{x}^{(i)}) = \mathcal{N}(\mathbf{x}^{(i)} + \frac{1}{2} \sigma_{\text{prop}}^2 \nabla_{\mathbf{x}} \log \Pr(\mathbf{x} | \mathbf{y}, \Sigma^{(\tau)})|_{\mathbf{x}^{(i)}}, \sigma_{\text{prop}}^2 I).$$

This proposal distribution is also straightforward to evaluate:

$$\begin{aligned} \nabla_{\mathbf{x}} \log \Pr(\mathbf{x} | \mathbf{y}, \Sigma^{(\tau)}) &= \nabla_{\mathbf{x}} \log \Pr(\mathbf{y} | \mathbf{x}, \Sigma^{(\tau)}) + \nabla_{\mathbf{x}} \log \Pr(\mathbf{x}) \\ &\approx H^\top \Sigma^{(\tau)-1} (\mathbf{y} - H\mathbf{x}) + \sigma_{\text{tr}}^{-2} (r(\mathbf{x}) - \mathbf{x}). \end{aligned}$$

METHODS AND RESULTS

Algorithm 1 Estimate latent signal \mathbf{x} and noise level Σ with the proposed methods AGEM and AGEM-ADMM. τ is the EM iteration number, initialized as 0. $\Sigma^{(1)}$ is initialized as $\sigma_{\text{tr}}^2 I$.

- Train a DAE with quadratic loss and noise $\eta \sim \mathcal{N}(0, \sigma_{\text{tr}}^2 I)$
- repeat $\tau \leftarrow \tau + 1$
- Initialization: If $\tau = 1$, $\mathbf{x}_\tau^{(1)} \leftarrow \mathbf{0}$, otherwise $\mathbf{x}_\tau^{(1)} \leftarrow \mathbf{x}_{\tau-1}^{(n_{\text{MH}})}$
- E-step: Draw n_{MH} samples $\{\mathbf{x}_\tau^{(i)}\}_{i=1}^{n_{\text{MH}}}$ with MALA, discard the first $\frac{1}{5}\tau$ samples as burn-in
- M-step: Use $\{\mathbf{x}_\tau^{(i)}\}_{i=n_{\text{MH}}/5}^{n_{\text{MH}}}$ to compute $\Sigma^{(\tau+1)}$
- until $\tau = n_{\text{EM}}$
- [AGEM] Compute $\hat{\mathbf{x}} \leftarrow \text{average of } \{\mathbf{x}_\tau^{(i)}\}_{i=n_{\text{MH}}/5}^{n_{\text{MH}}}$; return $(\hat{\mathbf{x}}, \Sigma^{(n_{\text{EM}})})$
- [AGEM-ADMM] Use ADMM and noise level $\Sigma^{(n_{\text{EM}})}$ to compute $\hat{\mathbf{x}}$; return $(\hat{\mathbf{x}}, \Sigma^{(n_{\text{EM}})})$

Results

Signal denoising. Consider 50-dimensional signals lying on a latent 2D manifold, and corrupted by isotropic Gaussian noise. The DAE is a multilayer perceptron with ReLU activations and 3 hidden layers, each containing 2000 neurons.

Table 1: Signal denoising, average RMSE of the test set. Standard deviations are in parentheses, estimated noise levels are in square brackets. Best performances are in bold. (All values are in 10^{-2}).

σ_n :	1.00		2.00		3.00		4.00	
	mean	std.	mean	std.	mean	std.	mean	std.
DAEP+NE [3]	0.73	(0.10)	0.98	(0.13)	1.16	(0.20)	1.31	(0.27)
ADMM+NE [35]	0.37	(0.28)	0.60	(0.36)	0.93	(0.55)	1.59	(3.49)
DMSP [4]	0.50	(0.22)	0.74	(0.29)	0.99	(0.45)	1.36	(0.95)
	[1.62]	(0.14)	[2.19]	(0.22)	[3.07]	(0.35)	[4.11]	(0.75)
AGEM	0.51	(0.15)	0.70	(0.25)	0.86	(0.39)	1.16	(0.64)
	[1.19]	(0.13)	[1.93]	(0.26)	[2.96]	(0.38)	[4.03]	(0.52)
AGEM-ADMM	0.33	(0.23)	0.57	(0.34)	0.91	(0.53)	1.43	(2.05)

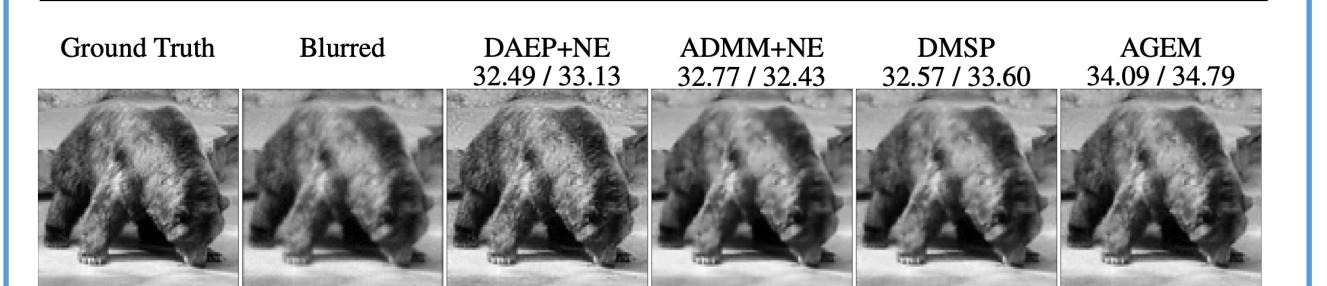
Our best method outperforms all baseline methods significantly statistically ($p < 0.01$), and its estimated σ_n are closer to true values than DMSP. We also compare with some analytic priors. Although these priors are specifically designed for image deconvolution, our generic approach outperforms them except for $\sigma_n = 0.04$, indicating that our trained DAE learns the distribution of natural images well, and DAE-based methods are indeed relevant in practice.

RESULTS

Image deblurring. We perform image deblurring with the STL-10 unlabeled dataset, which contains 105 colored 96×96 images. The DAE uses the full convolutional, residual architecture, where the input is added to the final layer's output. For testing, images are blurred using a 5×5 Gaussian filter with $\sigma = 0.6$.

Table 3: Average PSNR for image deblurring. Estimated noise levels are in square brackets.

σ_n :	0.01		0.02		0.03		0.04	
Method	mean	std.	mean	std.	mean	std.	mean	std.
DAEP+NE [3]	33.13	(1.39)	27.77	(0.89)	25.48	(0.70)	24.30	(0.61)
ADMM+NE [35]	32.43	(3.08)	29.48	(1.36)	27.87	(2.97)	25.78	(3.16)
DMSP [4]	33.60	(2.46)	30.89	(2.14)	28.93	(2.18)	27.40	(2.33)
	[0.017]	(1e-3)	[0.023]	(2e-3)	[0.027]	(3e-3)	[0.041]	(4e-3)
AGEM	34.79	(2.00)	31.42	(1.81)	29.47	(1.92)	28.00	(2.10)
	[0.014]	(1e-3)	[0.021]	(2e-3)	[0.030]	(3e-3)	[0.040]	(3e-3)
AGEM-ADMM	33.75	(2.70)	30.00	(3.20)	28.00	(2.88)	26.05	(3.51)



AGEM consistently outperforms all baseline methods significantly statistically ($p < 0.01$), and its estimated σ_n are closer to true values than DMSP. We also compare with some analytic priors. Although these priors are specifically designed for image deconvolution, our generic approach outperforms them except for $\sigma_n = 0.04$, indicating that our trained DAE learns the distribution of natural images well, and DAE-based methods are indeed relevant in practice.

