Quantifying Model Uncertainty in Inverse Problems via Bayesian Deep Gradient Descent



Introduction

- Deep learning (DL) image reconstruction techniques have remarkable results but lack estimates of uncertainty
- This is critical in sensitive domains such as medical imaging
- There are many "types or sources" of uncertainty but the most common in medical imaging are:
 - **Epistemic** uncertainty in the parameters (i.e., model uncertainty)
 - Aleatoric stochastic variability in the data
- **Our goal**: a DL reconstruction method that allows us to account for epistemic uncertainty in medical images

Deep Unrolled Optimisation

- Unrolled optimisation mimics iterative methods but
- 1. Executes only a finite number of iterations
- 2. Computes the "updates" using DNNs
- The iterates are computed as residual updates with a feasibility projection:

$$\mathbf{x}_k = \operatorname{ReLu}(\mathbf{x}_{k-1} + \delta \mathbf{x}_{k-1})$$

• The increments $\delta \mathbf{x}_{k-1}$ are computed as

$$\delta \mathbf{x}_{k-1} = f_{\varphi_k} \left(\nabla \mathcal{D} \left(\mathbf{y}, \mathbf{A} \mathbf{x}_{k-1} \right), \mathbf{x}_{k-1} \right) =: f_{\varphi_k} \left(\nabla \mathcal{D}, \mathbf{x}_k \right)$$

- $\varphi_k = (\phi_k, \theta_k)$ are (deterministic & probabilistic) parameters of (block) f_{φ_k}
- The entire iteration (**cascade**) consists of K sequential blocks
- $\mathbf{x}_{K} = (f_{\varphi_{K}} \circ f_{\varphi_{K-1}} \circ \cdots \circ f_{\varphi_{1}}) (\nabla \mathcal{D}, \mathbf{x}_{0}) := f_{\Phi_{K}}(\nabla \mathcal{D}, \mathbf{x}_{0}), \text{ with } \Phi_{K}$

Bayesian Deep Gradient Descent

- Each block (network) of the cascade consists of two parts • Deterministic layers with weights ϕ_k
- A final Bayesian layer with (random) weights θ_k
- For θ_k we learn the variational parameters ψ_k defining their distribution
- To estimate the posterior $p(\theta|X, Y)$ we use **variational inference** which uses an approximate, simple to compute, distribution q_{ψ}^{*}
- Moreover, we train the network **greedily**: provided previous k-1 blocks have been trained, in block k we use the family \mathcal{Q}_k

$$q_{\Psi_k}(\Theta_k) = q_{\Psi_{k-1}}^*(\Theta_{k-1})q_{\psi_k}(\theta_k|\Theta_{k-1}), \text{ with } q_{\psi_k}(\theta_k|\Theta_{k-1}) = \prod_{d=1}^D \mathcal{N}(\mu_{k,d}, \sigma_{k,d}^2),$$

where $\psi_k = \{(\mu_{k,\ell}, \sigma_{k,d}^2)\}_{d=1}^D$, and $q^*_{\Psi_{k-1}}$ is the distribution learnt for the previous k-1 blocks

The optimal distribution is computed by minimising the negative

$$q_{\Psi_k}^* \in \underset{q_{\Psi_k} \in \mathcal{Q}_k}{\operatorname{argmin}} \mathcal{L}_k(q_{\Psi_k}; X, Y) := -\int q_{\Psi_k}(\Theta_k) \log p(X|Y, \Theta_k) d\Theta_k + KI$$

• The prior is set recursively as

$$p(\Theta_k) = q_{\Psi_{k-1}}(\Theta_{k-1})p(\theta_k|\Theta_{k-1}), \text{ where } p(\theta_k|\Theta_{k-1}) = \mathcal{N}$$

BDGD takes the likelihood as:

$$p(\mathbf{x}|\mathbf{y},\Theta_k) = \mathcal{N}(f_{\Theta_k}(\nabla \mathcal{D},\mathbf{x}_0),\sigma_k^2 I)$$

Riccardo Barbano¹ Simon Arridge¹ Chen Zhang¹ Bangti Jin¹

¹University College London

Estimating Predictive Uncertainty

The optimal approximate posterior distribution is given by

$$q_{\Psi_K}^*(\Theta_K) = q_{\Psi_1}^*(\theta_1) \prod_{k=2}^K q_k$$

• We use Monte Carlo (MC) estimators to estimate the statistics of of the approximate predictive posterior defined as:

$$q^*(\mathbf{x}|\mathbf{y}) = \int p(\mathbf{x}|\mathbf{y}, \Theta_K) q_{\mathbf{x}}^*$$

• We compute the expected image (mean) and use $T \ge 1$ Monte Carlo samples:

$$\hat{\Xi}[\mathbf{x}] := \frac{1}{T} \sum_{t=1}^{T} f_{\hat{\Theta}_{K}^{t}}(\nabla$$

• We summarise predictive uncertainty as the (entry-wise) predictive variance $Var[\mathbf{x}]$ at the K^{th} step:

$$\operatorname{Var}[\mathbf{x}] = \operatorname{Var}_{q_{\Psi_{K}}(\Theta_{K})}[\mathbb{E}(\mathbf{x}|\mathbf{y},\Theta_{K})] + \mathbb{E}_{1}$$

$$\approx \sigma_K^2 + \frac{1}{T} \sum_{t=1}^{T} f_{\Theta_K^t} (\nabla \mathcal{D}, \mathbf{x}_0)^2 - \underbrace{\sum_{t=1}^{T} f_{\Theta_K^t} (\nabla \mathcal{D}, \mathbf{x}_0)^2}_{\text{oping}}$$

BDGD Framework Diagram



Practicalities in Training

Training

Input: # reconstruction steps K, dataset \mathcal{D} , initial guess $x_0^{(i)}$, batch-size M 1 for $k \leftarrow 1$ to K do

- Construct network's input:
- $\mathcal{D}_{k-1} = \{x_{k-1}^{(i)}, \nabla D(y^{(i)}, Ax_{k-1}^{(i)})\}_{i=1}^{N}$
- Train the k-th network $f_{\phi_k,\theta_k}(\nabla D(y^{(i)},Ax^{(i)}_{k-1}),x^{(i)}_{k-1})$: // stochastic mini-batch optimisation
- $\psi_k^*, \phi_k^* \leftarrow \arg\min_{q \in \mathcal{Q}_k, \phi_k} \left\{ \hat{\mathcal{L}}(\phi_k, q) = -\frac{N}{M} \sum_i^M \mathbb{E}_{\hat{\Theta}_k \sim q_{\Phi_k}(\Theta_k)} \left[\log p_{\Phi_k}(x^{(i)} | y^{(i)}, \hat{\Theta}_k) \right] + \mathrm{KL}(q_{\Phi_k}(\Theta_k) | | p_{\Phi_k}(\Theta_k)) \right\}$
- // update with $\hat{ heta}_k \sim q^*_{\Phi_k}(heta_k|\Theta_{k-1})$
- $x_{k}^{(i)} \leftarrow f_{\phi_{k},\hat{\theta}_{k}}(\nabla D(y^{(i)}, Ax_{k-1}^{(i)}), x_{k-1}^{(i)})$

Output: approximate posterior at each reconstruction step Inference

- **Input:** # reconstruction steps K, observation y, initial guess x_0 , trained parameters (Φ_K, Ψ_K), # samples T
- 9 for $t \leftarrow 1$ to T do //with $\hat{\Theta}_{K}^{(t)} \sim q_{\Phi_{K}}^{*}(\Theta_{K})$ 10
- 11
- **12** Evaluate $\hat{\mathbb{E}}[x]$ and $\widehat{\operatorname{Var}}[x]$ with $\{x_K^{(t)}\}_{t=1}^T$

Output: $\hat{\mathbb{E}}[x]$ and $\widehat{\operatorname{Var}}[x]$

s
$$\boldsymbol{\varphi}_k := (\varphi_1, \dots, \varphi_k)$$

 $L(q_{\Psi_k}(\Theta_k) \| p(\Theta_k))$

 $\mathcal{N}(0,I).$

 $q_{\Psi_k}^*(\theta_k|\Theta_{k-1})$

- $q_{\Psi_K}^*(\Theta_K) \mathrm{d}\Theta_K$
- $abla \mathcal{D}, \mathbf{x}_0)$

 $\mathbb{E}_{q_{\Psi_K}(\Theta_K)}[\operatorname{Var}(\mathbf{x}|\mathbf{y},\Theta_K)]$

 $rac{\mathbf{T}}{T}\sum f_{\Theta_K^t}(
abla \mathcal{D}, \mathbf{x}_0)$

Sparse view CT, with respect to the number of directions (dirs), and limited angle CT, with respect to the available range of angles. We report the mean PSNR over the ellipses dataset as well as the PSNR for the Shepp-Logan.

Table 1:Sparse View CT (30 dirs)			Table 2:Limited View CT $[0, 2\pi/3)$		
Methods	Ellipses Phantoms	SL Phantom	Methods	Ellipses Phantoms	SL Phantom
FBP	25.5264	18.4667	FBP	18.5958	17.1085
TV	35.1587	37.2162	TV	32.9134	29.2113
LPD	44.5122 ± 0.4911	44.0472 ± 0.4187	LPD [1]	40.7578 ± 0.3050	33.8427 ± 1.2380
DGD	43.2577 ± 0.4183	44.6913 ± 0.6644	DGD [2]	42.6994 ± 0.4243	42.8905 ± 0.5883
BDGD - MFVI	44.6642 ± 0.4637	47.2946 ± 0.5778	BDGD - MFVI	44.0297 ± 0.4698	45.5140 ± 0.8261
BDGD - MCDO	43.2126 ± 0.1285	45.1725 ± 0.4461	BDGD - MCDO	41.5367 ± 0.3884	41.4397 ± 0.6299



Fig. 1:Mean estimates and epistemic uncertainty maps by BDGD-MFVI for different geometries: (Left) sparse view with 30 directions, (*Centre*) limited view $[0, \pi/3)$, (*Right*) limited view $[0, 2\pi/3)$.



Fig. 2:Out-of-distribution reconstruction for different geometries by BDGD-MFVI: (Left) sparse view with 30 directions, (*Centre*) limited view $[0, \pi/3)$, (*Right*) limited view $[0, 2\pi/3)$.

- [1] Jonas Adler and Ozan Öktem. Learned primal-dual reconstruction. IEEE Trans. Med. Imag., 37(6):1322--1332, 2018.
- Ourselin, and Simon Arridge. IEEE Trans. Med. Imag., 37(6):1382--1393, 2018.

Results

Bibliography

[2] Andreas Hauptmann, Felix Lucka, Marta Betcke, Nam Huynh, Jonas Adler, Ben Cox, Paul Beard, Sebastien

Model-based learning for accelerated, limited-view 3-d photoacoustic tomography.