

Streamlining Prediction in Bayesian Deep Learning



Rui Li
Aalto University



Marcus Klasson
Aalto University



Arno Solin
Aalto University



Martin Trapp
Aalto University



Aalto University

International Conference on
Learning Representations

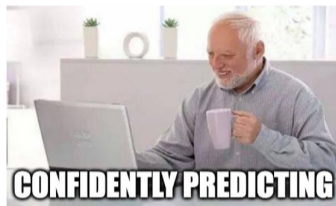


ICLR

Why Bayesian Deep Learning?

Deep learning is great, but it is not reliable:

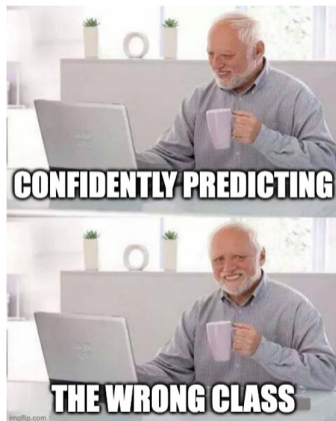
- ▶ **Sensitive** to perturbations
- ▶ **Don't know** when they don't know
- ▶ **Overconfident** predictions



Why Bayesian Deep Learning?

Deep learning is great, but it is not reliable:

- ▶ **Sensitive** to perturbations
- ▶ **Don't know** when they don't know
- ▶ **Overconfident** predictions

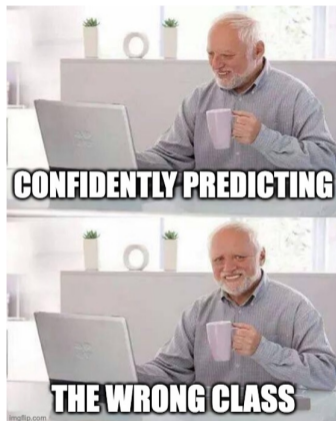


Why Bayesian Deep Learning?

Deep learning is great, but it is not reliable:

- ▶ **Sensitive** to perturbations
- ▶ **Don't know** when they don't know
- ▶ **Overconfident** predictions

💡 Bayesian deep learning aims to solve this.



Bayesian Deep Learning

1. Specify the prior: $p(\theta)$
2. Infer the posterior: $p(\theta | \mathcal{D})$
3. Make the prediction:

$$p(\mathbf{y}^* | \mathbf{x}^*) = \int p(\mathbf{y}^* | \mathbf{x}^*, \theta) p(\theta | \mathcal{D}) d\theta$$

Bayesian Deep Learning

1. Specify the prior: $p(\theta) = \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0)$
2. Infer the posterior: $p(\theta | \mathcal{D}) \approx \mathcal{N}(\mathbf{m}, \mathbf{S})$
3. Make the prediction:

$$p(\mathbf{y}^* | \mathbf{x}^*) = \int p(\mathbf{y}^* | \mathbf{x}^*, \theta) p(\theta | \mathcal{D}) d\theta$$

Bayesian Deep Learning

1. Specify the prior: $p(\theta) = \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0)$
2. Infer the posterior: $p(\theta | \mathcal{D}) \approx \mathcal{N}(\mathbf{m}, \mathbf{S})$
3. Make the prediction:

$$p(\mathbf{y}^* | \mathbf{x}^*) = \int p(\mathbf{y}^* | \mathbf{x}^*, \theta) p(\theta | \mathcal{D}) d\theta$$
$$\approx \underbrace{\frac{1}{S} \sum_{s=1}^S p(\mathbf{y}^* | \mathbf{x}^*, \theta^{(s)})}_{\text{(Challenge 2)}}, \underbrace{\theta^{(s)} \sim p(\theta | \mathcal{D})}_{\text{(Challenge 1)}}$$

Bayesian Deep Learning

1. Specify the prior: $p(\theta) = \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0)$
2. Infer the posterior: $p(\theta | \mathcal{D}) \approx \mathcal{N}(\mathbf{m}, \mathbf{S})$
3. Make the prediction:

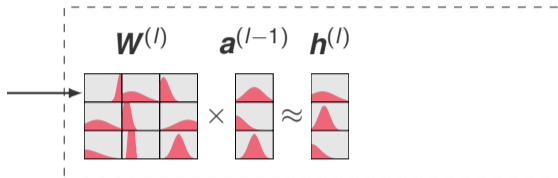
$$p(\mathbf{y}^* | \mathbf{x}^*) = \int p(\mathbf{y}^* | \mathbf{x}^*, \theta) p(\theta | \mathcal{D}) d\theta$$

~~$$\approx \frac{1}{S} \sum_{s=1}^S p(\mathbf{y}^* | \mathbf{x}^*, \theta^{(s)}), \quad \theta^{(s)} \sim p(\theta | \mathcal{D})$$~~

We approximate it with a Gaussian

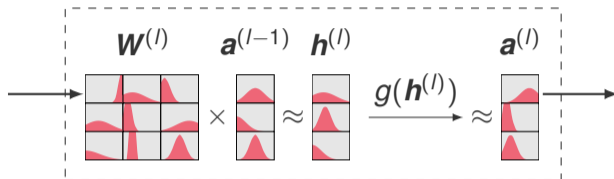


Streamlining Prediction in Bayesian Deep Learning



$$h^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)} \quad \text{approximate as Gaussian}$$

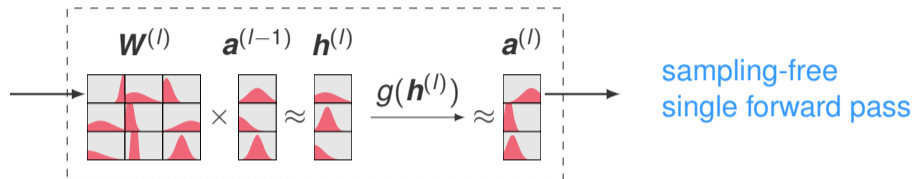
Streamlining Prediction in Bayesian Deep Learning



$$h^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)} \quad \text{approximate as Gaussian}$$

$$a^{(l)} = g(h^{(l)}) \quad \text{first order Taylor expansion}$$

Streamlining Prediction in Bayesian Deep Learning



$$h^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)} \quad \text{approximate as Gaussian}$$

$$a^{(l)} = g(h^{(l)}) \quad \text{first order Taylor expansion}$$

Efficient & Effective Uncertainty Quantification

Metrics	Methods	CIFAR-10	CIFAR-100	DTD	RESISC	IMAGENET-R
ACC \uparrow	LA Sampling	0.971 \pm 0.002	0.882 \pm 0.003	0.715 \pm 0.010	0.892 \pm 0.004	0.731 \pm 0.012
	LA GLM	0.976 \pm 0.002	0.879 \pm 0.003	0.718 \pm 0.010	0.891 \pm 0.004	0.739 \pm 0.012
	LA Ours	0.976 \pm 0.002	0.880 \pm 0.003	0.719 \pm 0.010	0.892 \pm 0.004	0.739 \pm 0.012
	MFVI Sampling	0.975 \pm 0.002	0.880 \pm 0.003	0.732 \pm 0.010	0.867 \pm 0.004	0.730 \pm 0.012
	MFVI Ours	0.975 \pm 0.002	0.880 \pm 0.003	0.734 \pm 0.010	0.867 \pm 0.004	0.728 \pm 0.012
NLPD \downarrow	LA Sampling	0.170 \pm 0.004	0.444 \pm 0.012	1.238 \pm 0.028	0.461 \pm 0.009	1.208 \pm 0.048
	LA GLM	0.092 \pm 0.007	0.459 \pm 0.012	1.197 \pm 0.029	0.385 \pm 0.010	1.180 \pm 0.047
	LA Ours	0.086 \pm 0.006	0.456 \pm 0.012	1.068 \pm 0.035	0.352 \pm 0.012	1.267 \pm 0.043
	MFVI Sampling	0.133 \pm 0.011	0.641 \pm 0.022	1.091 \pm 0.048	1.010 \pm 0.041	1.577 \pm 0.083
	MFVI Ours	0.088 \pm 0.006	0.468 \pm 0.013	1.007 \pm 0.035	0.617 \pm 0.019	1.234 \pm 0.052
ECE \downarrow	LA Sampling	0.006	0.022	0.197	0.129	0.070
	LA GLM	0.011	0.024	0.155	0.053	0.057
	LA Ours	0.008	0.027	0.040	0.016	0.132
	MFVI Sampling	0.015	0.070	0.075	0.079	0.118
	MFVI Ours	0.008	0.025	0.042	0.017	0.036

Efficient & Effective Uncertainty Quantification

Metrics	Methods	CIFAR-10	CIFAR-100	DTD	RESISC	IMAGENET-R
ACC \uparrow	LA Sampling	0.971 \pm 0.002	0.882 \pm 0.003	0.715 \pm 0.010	0.892 \pm 0.004	0.731 \pm 0.012
	LA GLM	0.976 \pm 0.002	0.879 \pm 0.003	0.718 \pm 0.010	0.891 \pm 0.004	0.739 \pm 0.012
	LA Ours	0.976 \pm 0.002	0.880 \pm 0.003	0.719 \pm 0.010	0.892 \pm 0.004	0.739 \pm 0.012
	MFVI Sampling	0.975 \pm 0.002	0.880 \pm 0.003	0.732 \pm 0.010	0.867 \pm 0.004	0.730 \pm 0.012
	MFVI Ours	0.975 \pm 0.002	0.880 \pm 0.003	0.734 \pm 0.010	0.867 \pm 0.004	0.728 \pm 0.012
NLPD \downarrow	LA Sampling	0.170 \pm 0.004	0.444 \pm 0.012	1.238 \pm 0.028	0.461 \pm 0.009	1.208 \pm 0.048
	LA GLM	0.092 \pm 0.007	0.459 \pm 0.012	1.197 \pm 0.029	0.385 \pm 0.010	1.180 \pm 0.047
	LA Ours	0.086 \pm 0.006	0.456 \pm 0.012	1.068 \pm 0.035	0.352 \pm 0.012	1.267 \pm 0.043
	MFVI Sampling	0.133 \pm 0.011	0.641 \pm 0.022	1.091 \pm 0.048	1.010 \pm 0.041	1.577 \pm 0.083
	MFVI Ours	0.088 \pm 0.006	0.468 \pm 0.013	1.007 \pm 0.035	0.617 \pm 0.019	1.234 \pm 0.052
ECE \downarrow	LA Sampling	0.006	0.022	0.197	0.129	0.070
	LA GLM	0.011	0.024	0.155	0.053	0.057
	LA Ours	0.008	0.027	0.040	0.016	0.132
	MFVI Sampling	0.015	0.070	0.075	0.079	0.118
	MFVI Ours	0.008	0.025	0.042	0.017	0.036

Efficient & Effective Uncertainty Quantification

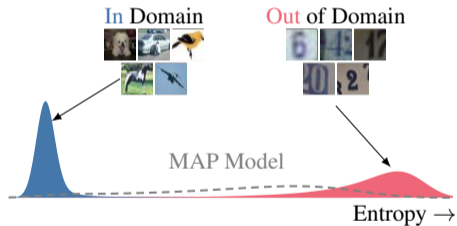
Metrics	Methods	CIFAR-10	CIFAR-100	DTD	RESISC	IMAGENET-R
ACC \uparrow	LA Sampling	0.971 \pm 0.002	0.882 \pm 0.003	0.715 \pm 0.010	0.892 \pm 0.004	0.731 \pm 0.012
	LA GLM	0.976 \pm 0.002	0.879 \pm 0.003	0.718 \pm 0.010	0.891 \pm 0.004	0.739 \pm 0.012
	LA Ours	0.976 \pm 0.002	0.880 \pm 0.003	0.719 \pm 0.010	0.892 \pm 0.004	0.739 \pm 0.012
	MFVI Sampling	0.975 \pm 0.002	0.880 \pm 0.003	0.732 \pm 0.010	0.867 \pm 0.004	0.730 \pm 0.012
	MFVI Ours	0.975 \pm 0.002	0.880 \pm 0.003	0.734 \pm 0.010	0.867 \pm 0.004	0.728 \pm 0.012
NLPD \downarrow	LA Sampling	0.170 \pm 0.004	0.444 \pm 0.012	1.238 \pm 0.028	0.461 \pm 0.009	1.208 \pm 0.048
	LA GLM	0.092 \pm 0.007	0.459 \pm 0.012	1.197 \pm 0.029	0.385 \pm 0.010	1.180 \pm 0.047
	LA Ours	0.086 \pm 0.006	0.456 \pm 0.012	1.068 \pm 0.035	0.352 \pm 0.012	1.267 \pm 0.043
	MFVI Sampling	0.133 \pm 0.011	0.641 \pm 0.022	1.091 \pm 0.048	1.010 \pm 0.041	1.577 \pm 0.083
	MFVI Ours	0.088 \pm 0.006	0.468 \pm 0.013	1.007 \pm 0.035	0.617 \pm 0.019	1.234 \pm 0.052
ECE \downarrow	LA Sampling	0.006	0.022	0.197	0.129	0.070
	LA GLM	0.011	0.024	0.155	0.053	0.057
	LA Ours	0.008	0.027	0.040	0.016	0.132
	MFVI Sampling	0.015	0.070	0.075	0.079	0.118
	MFVI Ours	0.008	0.025	0.042	0.017	0.036

Methods	AVG. RUNTIME (\pm STD) \downarrow
MAP	3.737 \pm 0.093
LA Sampling	190.806 \pm 0.137
LA GLM	17.191 \pm 0.734
MFVI Sampling	207.854 \pm 0.307
Ours (+ Cov)	14.728 \pm 0.144
Ours	4.350 \pm 0.079

On classification, we achieve better or on-par performance faster than baselines.

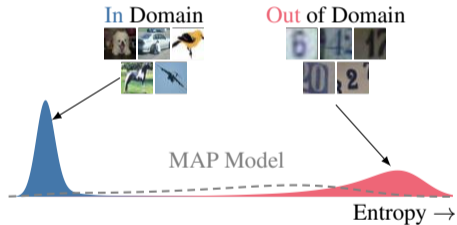
OOD detection and Input Sensitivity Analysis

Practical Outlier Detection

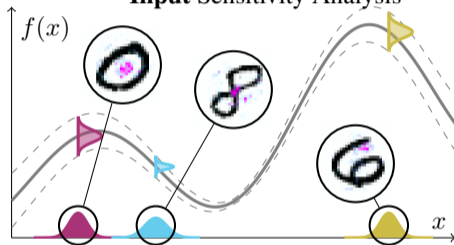


OOD detection and Input Sensitivity Analysis

Practical Outlier Detection



Input Sensitivity Analysis



Take away

Open-source library: <https://github.com/AaltoML/SUQ>

- ▶ **Goal:** Make **good** predictions **fast** in Bayesian neural networks.
- ▶ **Approach:** **Locally linearised** the neural network for a tractable posterior predictive distribution.
- ▶ **Result:** **Better** or on par performance with **faster** speed.