

Bayesian Deep Learning

Maurizio Filippone

maurizio.filippone@kaust.edu.sa

جامعة الملك عبد الله
للعلوم والتقنية
King Abdullah University of
Science and Technology



Statistics Program
King Abdullah University of Science and Technology

February 19, 2025



Decision-making is a critical step in several domains [**Norvig and Russell, 1995**]:

- ▶ Policy-making for the environment
- ▶ Healthcare
- ▶ Society
- ▶ ...



Decision-making is a critical step in several domains [**Norvig and Russell, 1995**]:

- ▶ Policy-making for the environment
- ▶ Healthcare
- ▶ Society
- ▶ ...

Decision Theory = Probabilistic reasoning + Utility theory



Decision-making is a critical step in several domains [**Norvig and Russell, 1995**]:

- ▶ Policy-making for the environment
- ▶ Healthcare
- ▶ Society
- ▶ ...

Decision Theory = Probabilistic reasoning + Utility theory

Is Deep Learning effective in assisting decision-making?

Over-confidence of Deep Learning Models - Online Meme



Image prediction: ping-pong ball
Confidence: 99.99%



Illustration: Dianna “Mick” McDougall, Photo: ResNeXtGuesser

Over-confidence of Deep Learning Models - Online Meme



Image prediction: pineapple
Confidence: 99.3%

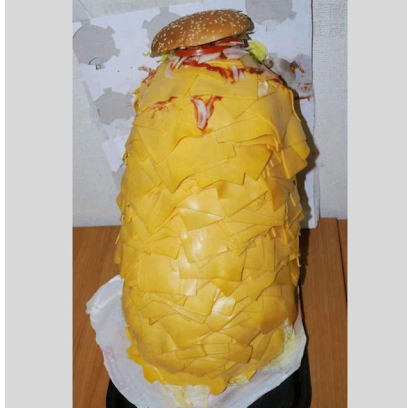
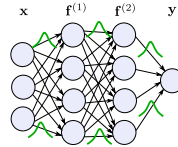


Illustration: Dianna “Mick” McDougall, Photo: ResNeXtGuesser



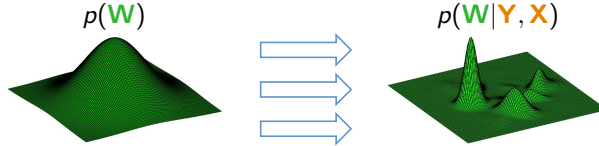
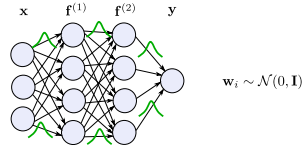
- Inputs : $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- Labels : $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$
- Weights : $\mathbf{W} = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}\}$



$$\mathbf{w}_i \sim \mathcal{N}(0, \mathbf{I})$$



- Inputs : $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- Labels : $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$
- Weights : $\mathbf{W} = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}\}$

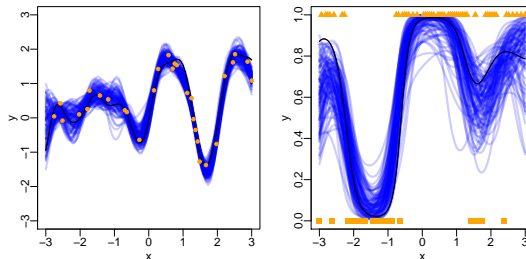


$$p(\mathbf{W}|\mathbf{Y}, \mathbf{X}) = \frac{p(\mathbf{Y}|\mathbf{X}, \mathbf{W})p(\mathbf{W})}{\int p(\mathbf{Y}|\mathbf{X}, \mathbf{W})p(\mathbf{W})d\mathbf{W}}$$



- Predictions consider an infinite number of parameter configurations (e.g., [Bishop, 2006])

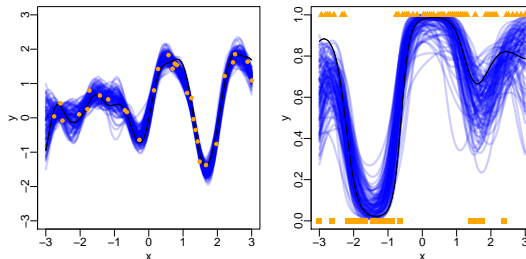
$$p(\mathbf{y}^* | \mathbf{x}^*, \mathbf{Y}, \mathbf{X}) = \int p(\mathbf{y}^* | \mathbf{x}^*, \mathbf{W}) p(\mathbf{W} | \mathbf{Y}, \mathbf{X}) d\mathbf{W}$$





- Predictions consider an infinite number of parameter configurations (e.g., [Bishop, 2006])

$$p(\mathbf{y}^* | \mathbf{x}^*, \mathbf{Y}, \mathbf{X}) = \int p(\mathbf{y}^* | \mathbf{x}^*, \mathbf{W}) p(\mathbf{W} | \mathbf{Y}, \mathbf{X}) d\mathbf{W}$$



Combining the flexibility of Deep Learning with sound uncertainty quantification!



- The normalization term in Bayes theorem can be used for model selection (e.g., [**Bishop, 2006; Gelman et al., 2013**])

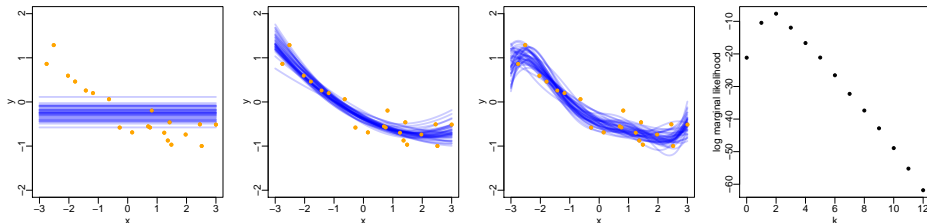
$$p(\mathbf{Y}|\mathbf{X}, \mathcal{H}) = \int p(\mathbf{Y}|\mathbf{X}, \mathbf{W}, \mathcal{H})p(\mathbf{W}|\mathcal{H})d\mathbf{W}$$



- The normalization term in Bayes theorem can be used for model selection (e.g., [Bishop, 2006; Gelman et al., 2013])

$$p(\mathbf{Y}|\mathbf{X}, \mathcal{H}) = \int p(\mathbf{Y}|\mathbf{X}, \mathbf{W}, \mathcal{H})p(\mathbf{W}|\mathcal{H})d\mathbf{W}$$

- Illustration of different hypotheses \mathcal{H} for a simple linear model:

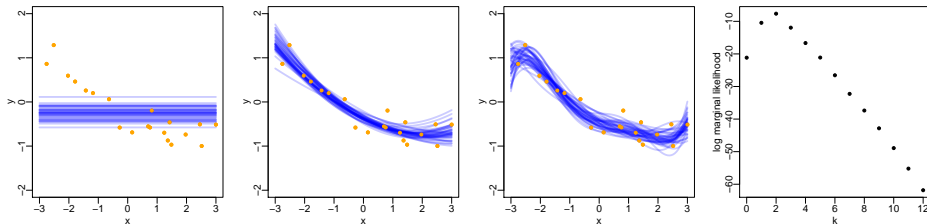




- The normalization term in Bayes theorem can be used for model selection (e.g., [Bishop, 2006; Gelman et al., 2013])

$$p(\mathbf{Y}|\mathbf{X}, \mathcal{H}) = \int p(\mathbf{Y}|\mathbf{X}, \mathbf{W}, \mathcal{H})p(\mathbf{W}|\mathcal{H})d\mathbf{W}$$

- Illustration of different hypotheses \mathcal{H} for a simple linear model:



Unfortunately, for Deep Learning models this is intractable!



- Leverage statistical connections [**Akaike, 1973**] between the model selection problem:

$$\arg \max_{\mathcal{H}} \{\log [p(\mathbf{Y}|\mathbf{X}, \mathcal{H})]\}$$

and the minimization of the following Kullback-Leibler divergence [**Tran et al., NeurIPS 2021**]:

$$\arg \min_{\mathcal{H}} \{\mathcal{D}_{\text{KL}}[\pi_{\text{TRUE}}(\mathbf{Y}|\mathbf{X}) \parallel p(\mathbf{Y}|\mathbf{X}, \mathcal{H})]\}$$

to obtain simple and tractable objectives for model selection.



- Leverage statistical connections [**Akaike, 1973**] between the model selection problem:

$$\arg \max_{\mathcal{H}} \{\log [p(\mathbf{Y}|\mathbf{X}, \mathcal{H})]\}$$

and the minimization of the following Kullback-Leibler divergence [**Tran et al., NeurIPS 2021**]:

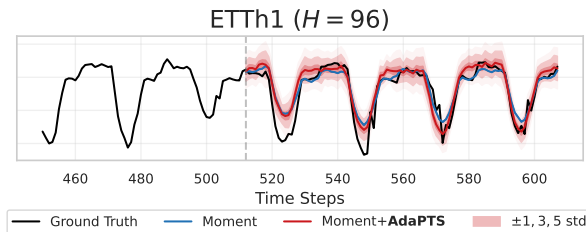
$$\arg \min_{\mathcal{H}} \{\mathcal{D}_{\text{KL}}[\pi_{\text{TRUE}}(\mathbf{Y}|\mathbf{X}) \parallel p(\mathbf{Y}|\mathbf{X}, \mathcal{H})]\}$$

to obtain simple and tractable objectives for model selection.

This can be leveraged for architecture search!

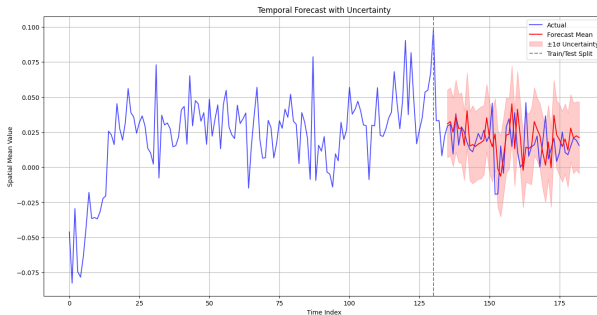


- ▶ Occam's razor [**Solomonoff, 1964**] for Deep Learning
- ▶ Efficient (compute and size) Deep Learning models
- ▶ Parsimonious probabilistic foundation models





► New foundation models for applications in Statistics





► Key collaborations in Deep Learning and Statistics worldwide

Position: Bayesian Deep Learning is Needed in the Age of Large-Scale AI

Theodore Papamarkou¹, Maria Skoularidou², Konstantina Palla³, Laurence Aitchison⁴, Julian Arbel⁵, David Dunson⁶, Maurizio Filippone⁷, Vincent Fortuin^{8,9,10}, Philipp Hennig¹¹, José Miguel Hernández-Lobato¹², Aliaksandr Hubin^{13,14}, Alexander Immer¹⁵, Theofanis Karaletsos¹⁶, Mohammad Emtiyaz Khan¹⁷, Agustinus Kristiadi¹⁸, Yingzhen Li¹⁹, Stephan Mandt²⁰, Christopher Nemeth²¹, Michael A. Osborne²², Tim G. J. Rudner²³, David Rüger²⁴, Yee Whye Teh^{25,26}, Max Welling²⁷, Andrew Gordon Wilson²⁸, Ruqi Zhang²⁹

Abstract

In the current landscape of deep learning research, there is a predominant emphasis on achieving high predictive accuracy in supervised tasks involving large image and language datasets. However, a broader perspective reveals a multitude of over-

uncertainty, active and continual learning, and scientific data, that demand attention. Bayesian deep learning (BDL) constitutes a promising avenue, offering advantages across these diverse settings. This paper posits that BDL can elevate the capabilities of deep learning. It revisits the strengths of BDL, acknowledges existing challenges, and

Spatial Statistics 60 (2024) 100825



Contents lists available at ScienceDirect

Spatial Statistics

journal homepage: www.elsevier.com/locate/spasta



Spatial Bayesian neural networks

Andrew Zammit-Mangion^{a,*}, Michael D. Kaminski^a, Ba-Hien Tran^b, Maurizio Filippone^c, Noel Cressie^a

^a School of Mathematics and Applied Statistics, University of Wollongong, Australia

^b Paris Research Centre, Huawei Technologies, France

^c Statistics Program, King Abdullah University of Science and Technology, Saudi Arabia

ARTICLE INFO

Keywords:

Gaussian process
Hamiltonian Monte Carlo

ABSTRACT

Statistical models for spatial processes play a central role in analyses of spatial data. Yet, it is the simple, interpretable, and well understood models that are routinely employed even though, as is pointed out in the introduction, more sophisticated models are available.

- Unique World-class AI and Statistics research environment!
- Vision 2030 drives methodological questions!



Thank you!

Questions?