# Machine Learning & Bayes - Synergies and Challenges -

Martin Trapp Aalto University

Learn Bayes Spring 2025 @ Karolinska Institutet

#### Outline



## Success of Modern Machine Learning

#### Generative tasks

#### Predictive tasks

#### Image Generation



#### Protein Prediction <sup>[1]</sup>

#### "Reasoning" tasks



Chain-of-Thought in LLMs <sup>[2]</sup>

#### Machine Learning in Healthcare



# Reminder on (modern) ML

... an incomplete and inaccurate picture of machine learning...

• The central objects in deep learning (modern ML) are artificial neural networks or neural network architectures (NNs).



6

• The central objects in deep learning (modern ML) are artificial neural networks or neural network architectures (NNs).



• The central objects in deep learning (modern ML) are artificial neural networks or neural network architectures (NNs).



- The central objects in deep learning (modern ML) are artificial neural networks or neural network architectures (NNs).
- NNs exploit function composition to define flexible non-linear function approximators with simple local operations.

$$f(\boldsymbol{x}) = f^{(L)} \circ \cdots \circ f^{(2)} \circ f^{(1)}(\boldsymbol{x})$$

- The central objects in deep learning (modern ML) are artificial neural networks or neural network architectures (NNs).
- NNs exploit function composition to define flexible non-linear function approximators with simple local operations.

$$f(\boldsymbol{x}) = f^{(L)} \circ \cdots \circ f^{(2)} \circ f^{(1)}(\boldsymbol{x})$$

• The network weights (parameters) are learned through backpropagation (chain rule) for a given loss function.

**O** PyTorch

In practice, modern NNs are complex systems:



Schematics of the Stable Diffusion 3.5 Model for Image Generation

https://huggingface.co/stabilityai/stable-diffusion-3.5-large

#### Machine Learning in Healthcare



#### Machine Learning in Healthcare

Many Applications are High-risk!



#### The Real World is Messy



ICCV 2023 Tutorial on The Many Faces of Reliability of Deep Learning for Real-World Deployment

#### Distribution Shifts

#### Various Sources of Noise



Koh, P W. et al. (2021) WILDS: A Benchmark of in-the-Wild Distribution Shifts. In *ICML*.



Sabour, S. et al. (2021) SpotlessSplats: Ignoring Distractors in 3D Gaussian Splatting. In *ICML*.



Out-of-domain data (Tesla)



Hallucinations (OpenAI Whisper)



#### N BIZ&IT CARS CULTURE GAMING HEALTH POLICY SCIENCE SECURITY SPACE TECH FORUM

#### TUNRELIABLE NARRATOR

#### Hospitals adopt errorprone AI transcription tools despite warnings

OpenAI's Whisper tool may add fake text to medical transcripts, investigation finds.

BENJ EDWARDS - 28 OCT 2024 20:23 | 🗩 169



-> Credit: Kobus Louw via Getty Images

## Expectations for ML Systems

Systems and models should be safe, trustworthy and reliable.

For this, we require that they are: *(incomplete list)* 

- Accurate in their predictions
- Robust to input perturbations (noise, adversarial attacks)
- "Know when they don't know" (recognise out-of-domain data)
- Act if they are uncertain (*e.g.*, active learning, reasoning)

#### However, ...

- Accurate in their predictions
- Sensitive to pertubations<sup>[1]</sup>
- Don't know when they don't know<sup>[2]</sup>
- Overconfident in their predictions<sup>[3]</sup>

Considered "unreliable"



Szegedy, C. et al. (2014). Intriguing properties of neural networks. In *ICLR*.
Nalisnick, E. et al. (2019). Do deep generative models know what they don't know? In *ICLR*.
Guo, C. et al. (2017). On calibration of modern neural networks. In *ICML*.



# Bayes for Machine Learning

#### **Functional Priors**



Meronen, L. et al. (2021). Periodic activation functions induce stationarity. In *NeurIPS*.

Functional Priors

#### Uncertainty Quant. & Overconfidence



Meronen, L. et al. (2021). Periodic activation functions induce stationarity. In *NeurIPS*.

Kristiadi, A. et al. (2020). Being Bayesian, Even Just a Bit, Fixes Overconfidence in ReLU Networks. In *ICML*.

Functional Priors

Uncertainty Quant. & Overconfidence

#### Active Learning



Meronen, L. et al. (2021). Periodic activation functions induce stationarity. In *NeurIPS*.



Kristiadi, A. et al. (2020). Being Bayesian, Even Just a Bit, Fixes Overconfidence in ReLU Networks. In *ICML*.



Baumann, A. et al. (2025). Post-hoc Probabilistic Vision-Language Models. ArXiv.

Functional Priors

#### Uncertainty Quant. & Overconfidence

#### Active Learning



Meronen, L. et al. (2021). Periodic activation functions induce stationarity. In NeurIPS.







Baumann, A. et al. (2025). Post-hoc Probabilistic Vision-Language Models. ArXiv.





## Uncertainty Quantification in ML



## Uncertainty Quantification in ML



## Uncertainty Quantification in ML









https://torch-uncertainty.github.io/



https://aleximmer.com/Laplace/



https://github.com/AaltoML/SUQ

(Deterministic) Neural Network



Bayesian Neural Network



(Deterministic) Neural Network

Bayesian Neural Network





Computations in Bayesian Learning are notoriously hard! \*in general.

Challenges in Bayesian Deep Learning:

- Hard to specify functional priors.
- Intractable high-dimensional multi-modal posterior. (billions of parameters)
- Models can be difficult to train from scratch. (high compute resource demand)

Computations in Bayesian Learning are notoriously hard! \*in general.

Challenges in Bayesian Deep Learning:

- Hard to specify functional priors.
- Intractable high-dimensional multi-modal posterior. (billions of parameters)
- Models can be difficult to train from scratch. (high compute resource demand)

#### Inference Methods in BDL



#### Inference Methods in BDL



#### Variational Inference: Big Picture

Recipe for approximating an intractable distribution  $p \in \mathcal{P}$ 

1. Define a tractable family of distributions Q

#### Variational Inference: Big Picture

Recipe for approximating an intractable distribution  $p \in \mathcal{P}$ 

- 1. Define a tractable family of distributions Q
- 2. Define a way to compute a "distance" between distributions

 $\mathrm{D}\left(p||\mathbf{q_1}\right) > \mathrm{D}\left(p||\mathbf{q_2}\right)$ 



#### Variational Inference: Big Picture

Recipe for approximating an intractable distribution  $p \in \mathcal{P}$ 

- 1. Define a tractable family of distributions Q
- 2. Define a way to compute a "distance" between distributions  $\mathbf{D}(\mathbf{u} \mid \mathbf{v}) = \mathbf{D}(\mathbf{u} \mid \mathbf{v})$

$$D\left(p||\boldsymbol{q}_{1}\right) > D\left(p||\boldsymbol{q}_{2}\right)$$

3. Search for best approximation

$$q^{\star} = \underset{q \in \mathcal{Q}}{\operatorname{arg\,min}\,\mathcal{D}}\left(p||q\right)$$



## Reminder on KL divergence

Let P and Q be two probability distributions with p.d.f. denoted as p and q, respectively.

Then their Kullback–Leibler divergence is given as:

$$D_{\mathrm{KL}}\left(P||Q\right) = \int_{x} p(x) \log\left(\frac{p(x)}{q(x)}\right) \mathrm{d}x$$

## Reminder on KL divergence

Let P and Q be two probability distributions with p.d.f. denoted as p and q, respectively.

Then their Kullback–Leibler divergence is given as:

$$D_{KL}(P||Q) = \int_{x} p(x) \log\left(\frac{p(x)}{q(x)}\right) dx$$

 $\mathcal{D}_{\mathrm{KL}}\left(P||Q\right) \geq 0 \qquad \text{ and } \qquad \mathcal{D}_{\mathrm{KL}}\left(P||Q\right) = 0 \iff P = Q$ 

Reverse KL  $D_{\mathrm{KL}}\left(q(\theta) || p(\theta \mid \boldsymbol{x}_{1:n})\right) = \mathbb{E}_{\theta \sim q}\left[\log\left(\frac{q(\theta)}{p(\theta \mid \boldsymbol{x}_{1:n})}\right)\right]$ 

$$D_{\mathrm{KL}}(q(\theta)||p(\theta \mid \boldsymbol{x}_{1:n})) = \mathbb{E}_{\theta \sim q} \left[ \log \left( \frac{q(\theta)}{p(\theta \mid \boldsymbol{x}_{1:n})} \right) \right]$$
$$= \mathbb{E}_{\theta \sim q} \left[ \log \left( \frac{q(\theta)}{\frac{p(\boldsymbol{x}_{1:n} \mid \theta) p(\theta)}{p(\boldsymbol{x}_{1:n})}} \right) \right]$$

$$D_{\mathrm{KL}}(q(\theta)||p(\theta \mid \boldsymbol{x}_{1:n})) = \mathbb{E}_{\theta \sim q} \left[ \log \left( \frac{q(\theta)}{p(\theta \mid \boldsymbol{x}_{1:n})} \right) \right]$$
$$= \mathbb{E}_{\theta \sim q} \left[ \log \left( \frac{q(\theta)}{\frac{p(\boldsymbol{x}_{1:n} \mid \theta) p(\theta)}{p(\boldsymbol{x}_{1:n})}} \right) \right]$$
$$= \mathbb{E}_{\theta \sim q} \left[ \log \left( \frac{q(\theta)}{p(\theta)} \right) + \log \left( \frac{p(\boldsymbol{x}_{1:n})}{p(\boldsymbol{x}_{1:n} \mid \theta)} \right) \right]$$

$$\begin{aligned} \mathbf{D}_{\mathrm{KL}}\left(q(\theta)||p(\theta \mid \boldsymbol{x}_{1:n})\right) &= \mathbb{E}_{\theta \sim q} \left[\log\left(\frac{q(\theta)}{p(\theta \mid \boldsymbol{x}_{1:n})}\right)\right] \\ &= \mathbb{E}_{\theta \sim q} \left[\log\left(\frac{q(\theta)}{\frac{p(\boldsymbol{x}_{1:n}\mid\theta) p(\theta)}{p(\boldsymbol{x}_{1:n})}}\right)\right] \\ &= \mathbb{E}_{\theta \sim q} \left[\log\left(\frac{q(\theta)}{p(\theta)}\right) + \log\left(\frac{p(\boldsymbol{x}_{1:n})}{p(\boldsymbol{x}_{1:n}\mid\theta)}\right)\right] \\ &= \mathbf{D}_{\mathrm{KL}}\left(q(\theta)||p(\theta)\right) - \mathbb{E}_{\theta \sim q} \left[\log p(\boldsymbol{x}_{1:n}\mid\theta)\right] + \log p(\boldsymbol{x}_{1:n}) \end{aligned}$$

$$D_{\mathrm{KL}}\left(q(\theta) || p(\theta \mid \boldsymbol{x}_{1:n})\right) = \mathbb{E}_{\theta \sim q}\left[\log\left(\frac{q(\theta)}{p(\theta \mid \boldsymbol{x}_{1:n})}\right)\right]$$

•

 $= \mathrm{D}_{\mathrm{KL}} \left( q(\theta) || p(\theta) \right) - \mathbb{E}_{\theta \sim q} \left[ \log p(\boldsymbol{x}_{1:n} \mid \theta) \right] + \log p(\boldsymbol{x}_{1:n})$ 

$$D_{\mathrm{KL}}\left(q(\theta) || p(\theta \mid \boldsymbol{x}_{1:n})\right) = \mathbb{E}_{\theta \sim q}\left[\log\left(\frac{q(\theta)}{p(\theta \mid \boldsymbol{x}_{1:n})}\right)\right]$$

•

 $= \mathrm{D}_{\mathrm{KL}} \left( q(\theta) || p(\theta) \right) - \mathbb{E}_{\theta \sim q} \left[ \log p(\boldsymbol{x}_{1:n} \mid \theta) \right] + \log p(\boldsymbol{x}_{1:n})$ 

 $\log p(\boldsymbol{x}_{1:n}) = -\mathrm{D}_{\mathrm{KL}} \left( q(\theta) || p(\theta) \right) + \mathbb{E}_{\theta \sim q} \left[ \log p(\boldsymbol{x}_{1:n} \mid \theta) \right] + \underbrace{\mathrm{D}_{\mathrm{KL}} \left( q(\theta) || p(\theta \mid \boldsymbol{x}_{1:n}) \right)}_{>0}$ 

$$D_{\mathrm{KL}}\left(q(\theta) || p(\theta \mid \boldsymbol{x}_{1:n})\right) = \mathbb{E}_{\theta \sim q}\left[\log\left(\frac{q(\theta)}{p(\theta \mid \boldsymbol{x}_{1:n})}\right)\right]$$

 $= \mathrm{D}_{\mathrm{KL}} \left( q(\theta) || p(\theta) \right) - \mathbb{E}_{\theta \sim q} \left[ \log p(\boldsymbol{x}_{1:n} \mid \theta) \right] + \log p(\boldsymbol{x}_{1:n})$ 

 $\log p(\boldsymbol{x}_{1:n}) = -\mathrm{D}_{\mathrm{KL}} \left( q(\theta) || p(\theta) \right) + \mathbb{E}_{\theta \sim q} \left[ \log p(\boldsymbol{x}_{1:n} \mid \theta) \right] + \underbrace{\mathrm{D}_{\mathrm{KL}} \left( q(\theta) || p(\theta \mid \boldsymbol{x}_{1:n}) \right)}_{\geq 0}$ 

 $\geq -D_{\mathrm{KL}}\left(q(\theta)||p(\theta)\right) + \mathbb{E}_{\theta \sim q}\left[\log p(\boldsymbol{x}_{1:n} \mid \theta)\right] = \mathrm{ELBO} \text{ (Evidence Lower Bound)}$ 













https://lacerbi.github.io/blog/2024/vi-is-inference-is-optimization/

## Variational Inference: Takeaways

- VI turns inference into an optimisation problem.
- Can scale to large data sets, but is challenging to optimise
- Can exhibit high variance in stochastic gradients
- Tendency to underestimate the variance if reverse KL is used
- Recent work has shown that VI can be effective for large-scale models.

What if training a model "from scratch" is too expensive?



What if training a model "from scratch" is too expensive?







$$\underbrace{\log p(\theta) + \log p(\boldsymbol{x} \mid \theta)}_{\ell(\theta)} \approx \ell(\theta^*) + J_{\ell|_{\theta=\theta^*}}(\theta - \theta^*) - \frac{1}{2}H_{\ell|_{\theta=\theta^*}}(\theta - \theta^*)^2$$



$$\underbrace{\log p(\theta) + \log p(\boldsymbol{x} \mid \theta)}_{\ell(\theta)} \approx \ell(\theta^*) + J_{\ell|_{\theta=\theta^*}}(\theta - \theta^*) - \frac{1}{2}H_{\ell|_{\theta=\theta^*}}(\theta - \theta^*)^2$$

$$= \ell(\theta^*) - \frac{1}{2} \boldsymbol{H}_{\ell|_{\theta=\theta^*}} (\theta - \theta^*)^2$$
  
 
$$\propto \log \mathcal{N}(\theta^*, \boldsymbol{H}_{\ell|_{\theta=\theta^*}}^{-1})$$



Daxberger, E. et al. (2021) Laplace Redux - Effortless Bayesian Deep Learning. In *NeurIPS*. Roy, S. et al. (2022). Uncertainty-guided Source-free Domain Adaptation. In *ECCV*.

#### Linearised Laplace: Problem Setting



Antoran, J. (2024). Scalable Bayesian Inference in the Era of Deep Learning: From Gaussian Processes to Deep Neural Networks. PhD thesis.

#### Linearised Laplace: Problem Setting



Antoran, J. (2024). Scalable Bayesian Inference in the Era of Deep Learning: From Gaussian Processes to Deep Neural Networks. PhD thesis.

#### Linearised Laplace

The linearised Laplace approximates the predictive distribution:

$$f(x;\theta) \approx f(x;\theta_{\mathrm{MAP}}) + \nabla f(x;\theta) \mid_{\theta_{\mathrm{MAP}}}^{\top} (\theta - \theta_{\mathrm{MAP}}) = \hat{f}(x;\theta)$$

#### Linearised Laplace

The linearised Laplace approximates the predictive distribution:

$$f(x;\theta) \approx f(x;\theta_{\mathrm{MAP}}) + \nabla f(x;\theta) \mid_{\theta_{\mathrm{MAP}}}^{\top} (\theta - \theta_{\mathrm{MAP}}) = \hat{f}(x;\theta)$$

$$p(y^{\star} \mid x^{\star}, \boldsymbol{x}, \boldsymbol{y}) \approx \int \mathcal{N}(y^{\star} \mid \hat{f}(\boldsymbol{x}; \boldsymbol{\theta}), \sigma) \mathcal{N}(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\mathrm{MAP}}, \Sigma) \mathrm{d}\boldsymbol{\theta}$$

## Laplace Approximation: Takeaways

- Does not require retraining (post-hoc).
- Requires linearisation in the case of deep learning models.
- Crude and only local approximation.
- Estimation of the Hessian can be difficult.

• Crude but useful (easy to use) "tool" for Bayesian inference in deep learning.

# Machine Learning for Bayes

## **Amortized Inference**

- NNs are powerful function approximators and generators (e.g., diffusion models)
- How can we use NNs for approximate Bayesian inference?
- $\rightarrow$  Emerging research on amortized inference with NNs (condition and predict)



(a) Prior v samples

(b) PriorGuide v samples

(c) PriorGuide v retrained

Chang, P. et al. (2025). Inference-Time Prior Adaptation in Simulation-Based Inference via Guided Diffusion Models. ICLR workshop.

# Recap

- Deep learning methods can be problematic in high-risk domains.
- The Bayesian approach to deep learning can help reduce overconfidence, quantify uncertainties, and utilise uncertainties in decision-making.
- Prior specification is challenging and an open question.
- Performing computations is challenging, but promising avenues exist.
- Machine learning can help scaling Bayesian inference, e.g., through amortisation.

#### Thanks for your attention!

