

Introduction to Bayesian Deep Learning

Andreas Makris, 2nd year PhD student
Lancaster University, ProbAI Hub

Based on new Bayesian Deep Learning book!

Term 2 Schedule (So far)

Schedule

Title	Location	Date	Time	Speaker
Intro to BDL	PSC LT	14 Jan 2026	2 pm - 3 pm	Andreas
Laplace Methods	PSC LT	21 Jan 2026	2 pm - 3 pm	Lanya
Robust & Conjugate Spatio-Temporal GPs	PSC LT	28 Jan 2026	2 pm - 3 pm	William
TBD	PSC LT	4 Feb 2026	2 pm - 3 pm	Paul
Reviewing BNNs	PSC LT	11 Feb 2026	2 pm - 3 pm	Jixiang
GP Connections to BDL	PSC LT	27 Feb 2026	2 pm - 3 pm	Jonas
TBD	PSC LT	4 Mar 2026	2 pm - 3 pm	Connie
TBD	PSC LT	11 Mar 2026	2 pm - 3 pm	Chris
Probabilistic Numerics	PSC LT	18 Mar 2026	2 pm - 3 pm	Rui

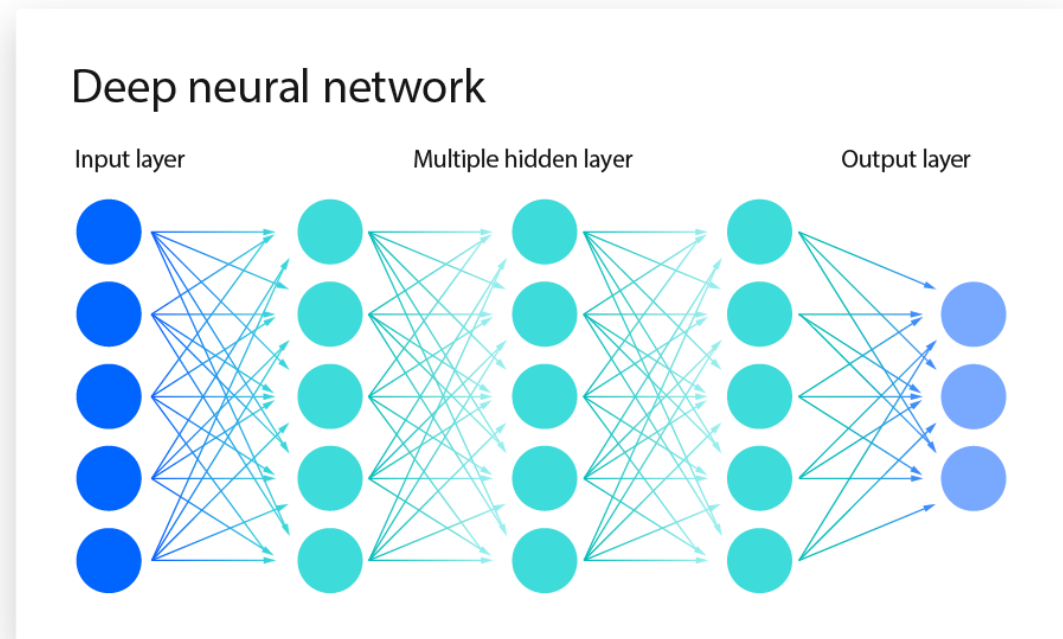
+ potentially an active learning talk

+ potentially a last layer inference talk

Why are we even here?

- **Deep Learning works extremely well. What is it missing?**
- Most neural networks give a point prediction.
- They are often extremely confident, even when they are wrong.
- They do not naturally tell us when they should not be trusted.

This is not a minor issue; it becomes **critical whenever predictions are used to make important decisions** (e.g. medical cases).



When is a point prediction not enough?

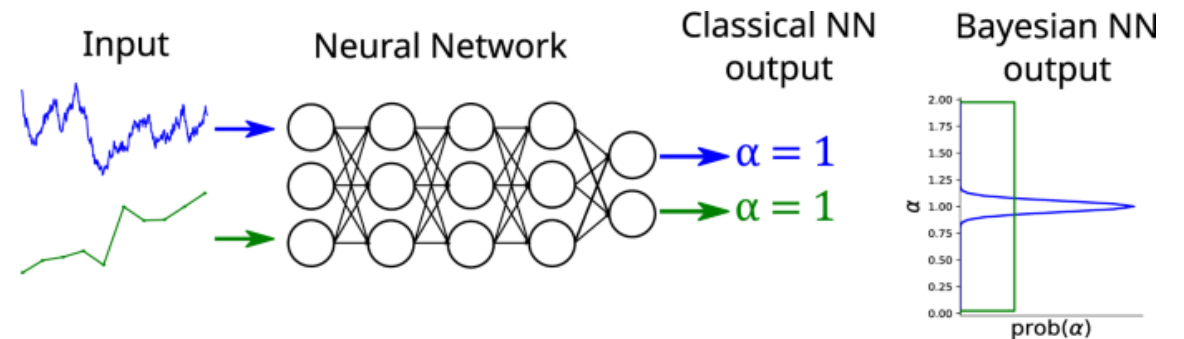
- Medical diagnosis.
- Autonomous driving.
- LLMs (hallucinations).
- Spatio temporal forecasting (uncertainty propagates over time).

In all of these cases (and many more), we don't just want a prediction. We want to know **how certain that prediction is**.



Bayesian Deep Learning (BDL)

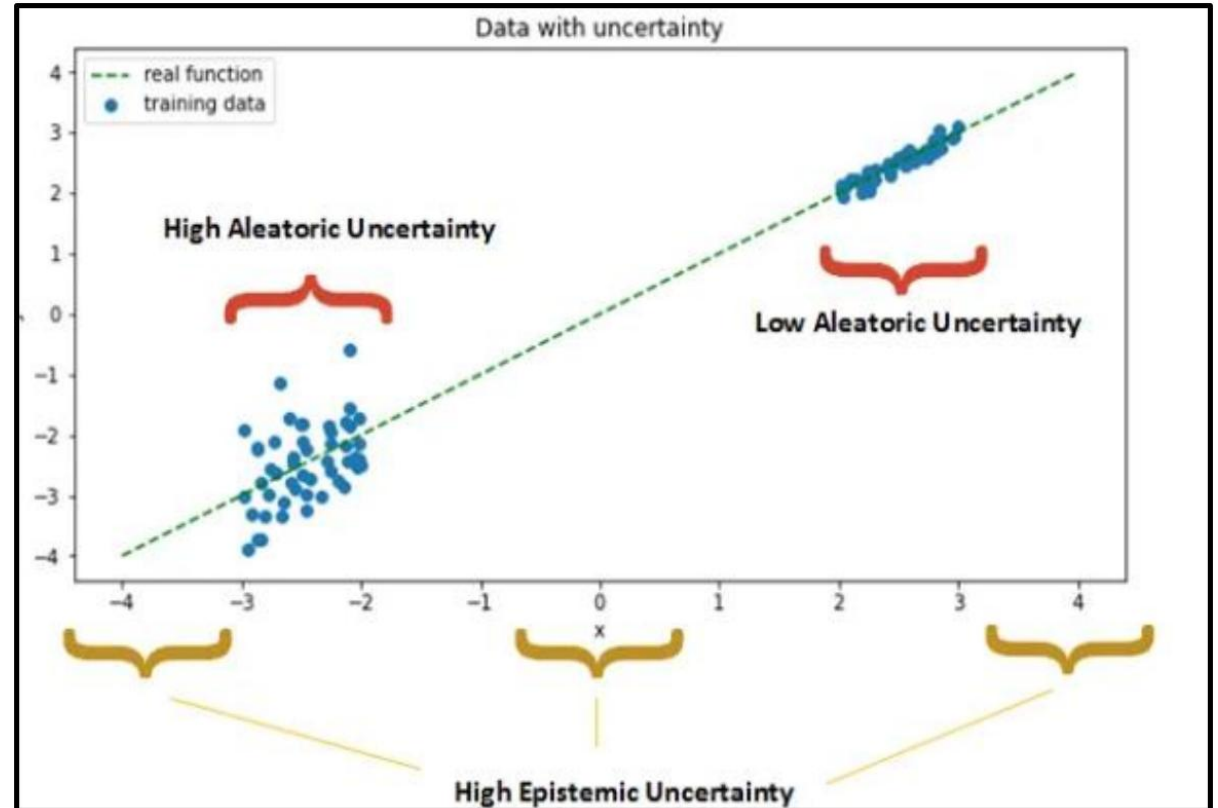
- “Bayesian deep learning **integrates deep neural networks with Bayesian inference**. Unlike traditional deep learning, which produces deterministic predictions, Bayesian deep learning **treats model parameters as random variables** with associated probability distributions. This **allows the model to capture and quantify uncertainty** inherent in data, leading to more robust and reliable predictions.”



Type of Uncertainty

- Epistemic Uncertainty:
 - Due to limited data.
 - Reducible with more data.
- Aleatoric Uncertainty:
 - Inherent noise in the data.
 - Irreducible.
 - Measured by e.g. output of neural network are parameters of a Gaussian.

Bayesian deep learning is primarily concerned with epistemic uncertainty. This is exactly the uncertainty that classical deep learning ignores.



Neural Network Parameters

We have some data $\mathcal{D} = \{X, Y\}$. We want to find the posterior distribution of our neural network parameters θ .

$$p(\theta|\mathcal{D}) = \frac{p(Y|\theta, X)p(\theta)}{p(\mathcal{D})}$$

Our ultimate goal is estimating the posterior predictive distribution for a new input x^* .

$$p(y^*|x^*, \mathcal{D})$$

$$\boxed{P(A|B)}_{\text{posterior}} = \boxed{P(A)}_{\text{prior}} \times \frac{\boxed{P(B|A)}_{\text{likelihood}}}{\boxed{P(B)}_{\text{marginal}}}$$

Estimating the Posterior

Standard deep learning only finds a **MAP** of the posterior. To better estimate the posterior, we can use:

1. Sampling-based approximations (SG-MCMC).
2. Laplace approximations.
3. Variational Inference.
4. Ensemble methods.

Due to the scale of the models, **no method gives exact samples.**

- MCMC is exact only asymptotically (computationally infeasible for deep learning).
- SG-MCMC is an approximation of MCMC.
- Laplace approximations, VI are considered approximate Bayesian inference.
- Ensemble methods are “non-Bayesian”, but in practice they provide better approximations than traditional approximate Bayesian methods.

SG-MCMC

- We usually can't use simple sampling techniques (like rejection sampling, MCMC) due to the high dimensionality of the parameters and the size of the dataset.
- We can use Stochastic Gradient MCMC (SG-MCMC), a class of methods that **uses stochastic gradients to make MCMC suitable for deep learning models**.
- Compared to mini-batch GD, an additional noise term is introduced.

Algorithm 1: SGLD

Input: $\theta_0, \{h_0, \dots, h_K\}$.

for $k \in 1, \dots, K$ **do**

 Draw $\mathcal{S}_n \subset \{1, \dots, N\}$ without replacement

 Estimate $\hat{\nabla}U(\theta)^{(n)}$ using (4)

 Draw $\xi_k \sim N(0, h_k I)$

 Update $\theta_{k+1} \leftarrow \theta_k - \frac{h_k}{2} \hat{\nabla}U(\theta_k)^{(n)} + \xi_k$

end

$$\hat{\nabla}U(\theta)^{(n)} = \frac{N}{n} \sum_{i \in \mathcal{S}_n} \nabla U_i(\theta), \quad (4)$$

Laplace Approximations

Using sampling methods has huge memory requirements.

- Estimate the posterior by a Gaussian distribution.

$$p(\theta|\mathcal{D}) \approx \mathcal{N}(\theta_*, \Psi^{-1}(\theta_*))$$

- If we have a trained deep learning model, we can just use its parameters as θ_* .
- We can compute the Hessian using autodiff **for small models**.
- For large models, we need to approximate it with e.g. a diagonal matrix (next week).
- The memory and compute costs, and quality of the approximation of the Hessian are controlled by the user.

- We can do the Laplace approximation post hoc (if we have access to the training data)!
- Simple way to fix some nontrivial pathologies of standard deep learning.

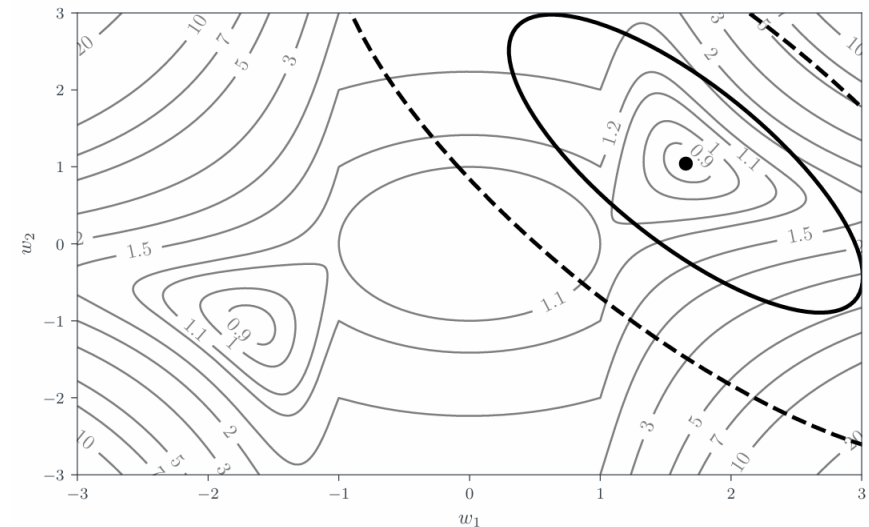


Figure 7.1 A linearized Laplace approximation on the most basic feed-forward perceptron $f(x, \theta) = \theta_2 \cdot \text{ReLU}(\theta_1 x + 1)$ for the dataset $X = (1, -1), Y = (-1, -1)$, with regularizer $\Omega(\theta) = 0.2 \cdot \|\theta\|^2$ and square loss. The gray contour lines denote the Energy function. The black dot represents the minimum θ_* ; the thick black ellipses denote one and two standard deviations of the Laplace approximation.

Variational Inference

- VI reframes the Bayesian inference as an **optimization problem**.
- We introduce a variational family of distributions (user defined hyperparameter, simplest and most common case is a normal distribution) $q(\theta; \lambda)$ and find the “best” choice of the parameters λ .

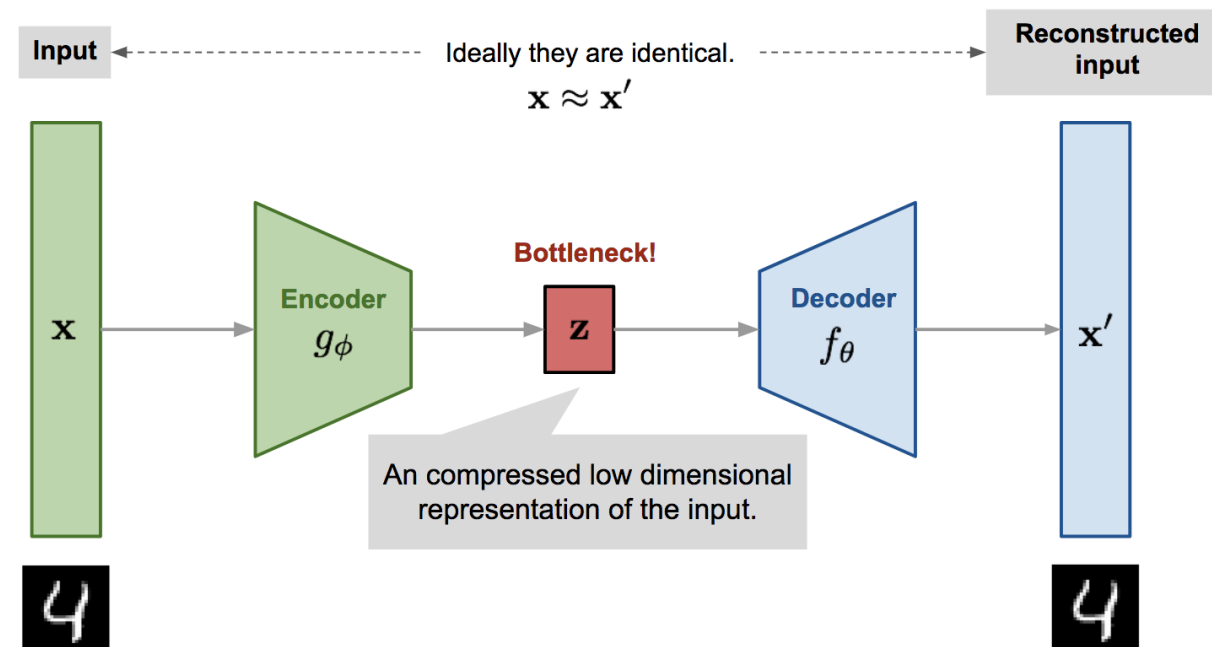
$$\lambda^* = \operatorname{argmin}_{\lambda} D(q(\theta; \lambda) || p(\theta | \mathcal{D}))$$

D is a divergence (similar to a distance but is not symmetric). An example is the KL-Divergence.

A simple choice of q is:

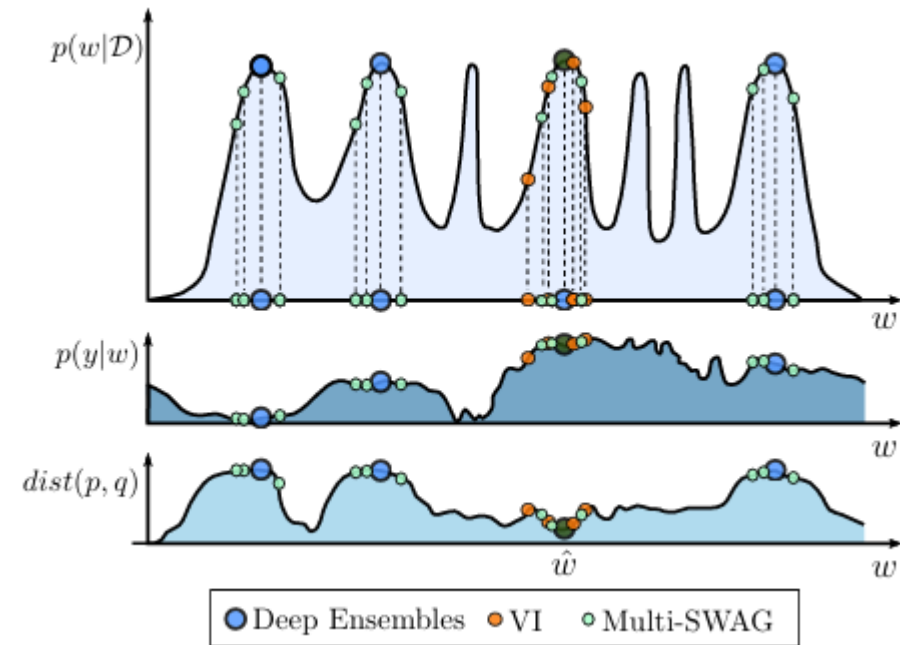
$$q(\theta; \lambda) = \prod \mathcal{N}(\theta_i | \mu_i, \sigma_i^2)$$

VAEs (not BDL). They are trained with VI.



Deep Ensembles

- Retrain a neural network with different random initializations and ensemble the results.
- Typically works better than standard approximate Bayesian methods.
- If we use K ensembles, then we have an additional K -fold training and testing cost.
- *What if instead we use the K -fold increase in computation to train a larger neural network?*
 - It seems that often, the larger neural network will yield better accuracy.



Last Layer Inference

- Having a posterior distribution over each parameter is very expensive, especially for large networks.
- Instead treat all parameters as being deterministic, **except the ones in the last layer**.
- This can be interpreted as performing Bayesian linear regression with covariates the output of a neural network (all but the last layer).
- So much cheaper during inference!!! We only pass through the model once (except the last layer) and perform multiple passes only through the last layer!

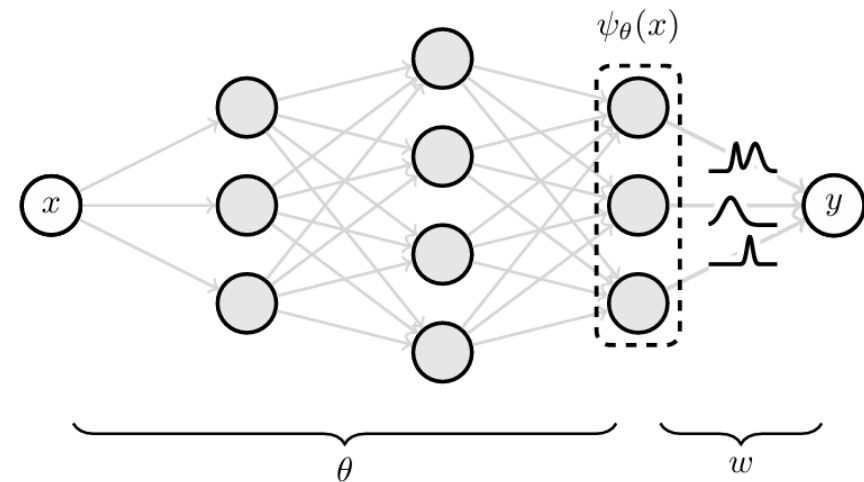
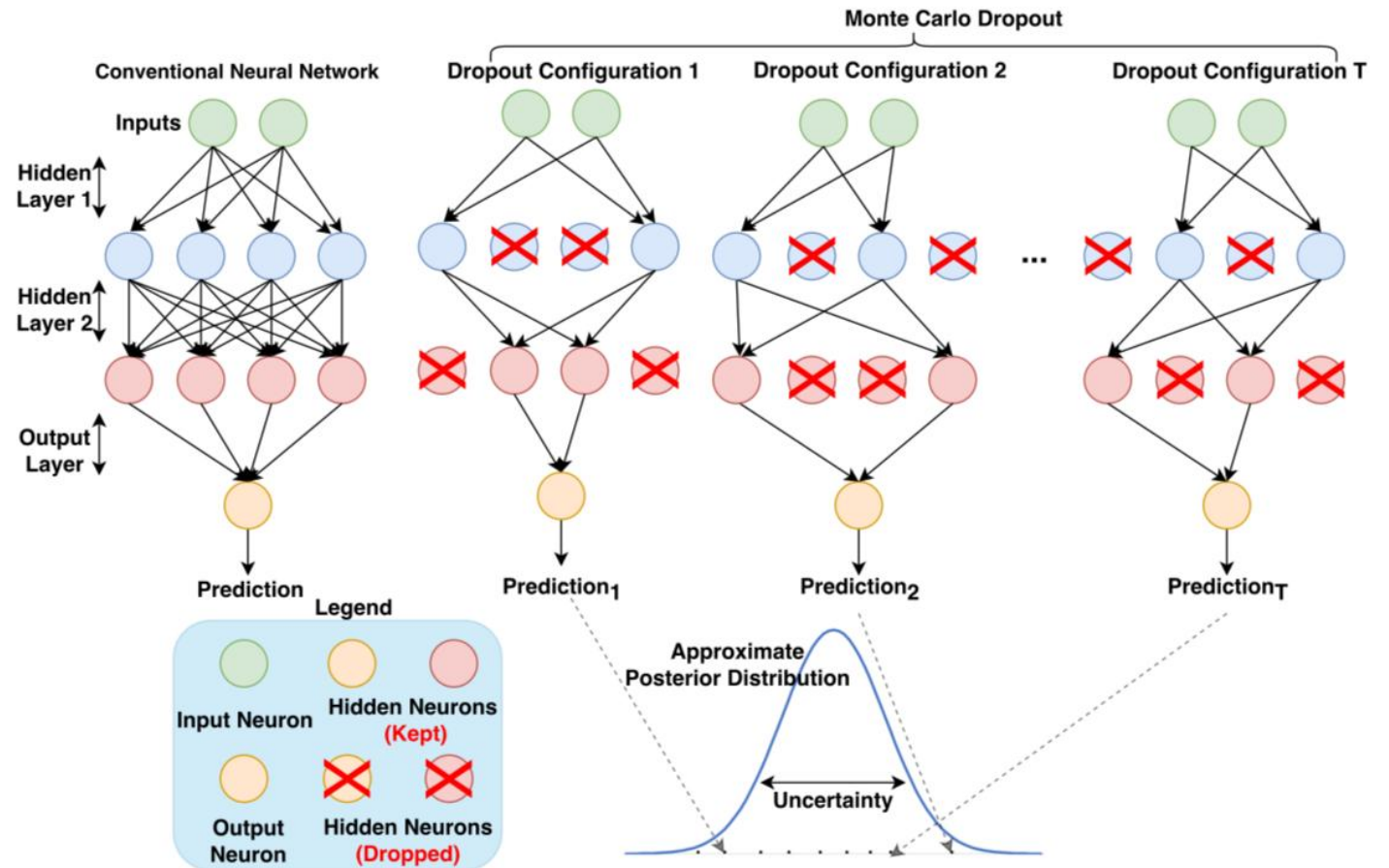


Figure 39.1 In last-layer inference, the parameters $\theta \in \mathbb{R}^{P-H_L-1}$ up to and including the penultimate layer are treated as deterministic. Only the posterior $p(w \mid \mathcal{D})$ for the last-layer weights is computed using the deep feature projection ψ_θ . Figure inspired by **blundell2015weight**.

MC Dropout

- Train a model with dropout (or use a pretrained one).
- Perform multiple passes during inference **with dropout**.



Gaussian Processes

- GPs are not deep learning (yet).
- GPs extend the notion of a Normal distribution to functions (infinite dimension).
- In practice, we often set $m(x) = 0$.
- The kernel is a measure of similarity, e.g. RBF kernel.

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\ell^2}\right)$$

Definition 23.3 (Gaussian Process [Rasmussen2006Gaussian]). Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be a real-valued function defined on an input space \mathcal{X} . We say that f is a Gaussian Process, denoted

$$f \sim \mathcal{GP}(m(\cdot), k(\cdot, \cdot)),$$

if for any finite collection of inputs $\mathbf{x} = [x_1, \dots, x_n]^\top$, the random vector of function evaluations

$$\mathbf{f} = [f(x_1), \dots, f(x_n)]^\top$$

follows a multivariate normal distribution

$$\mathbf{f} \sim \mathcal{N}(\boldsymbol{\mu}, K_{\mathbf{xx}}),$$

with mean $\boldsymbol{\mu}_i = m(x_i)$ and covariance $[K_{\mathbf{xx}}]_{ij} = k(x_i, x_j)$.

Prior

“The user’s choice of kernel amounts to specifying a prior over function behaviour and determining what kinds of variation are expected in unseen regions of the input space.”

Gaussian Processes for Regression

Given training data $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^P$, we assume that observations are noisy evaluations of an underlying latent function f :

$$y_i = f(x_i) + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2),$$

where σ^2 denotes the noise variance.

Let $\mathbf{y} = [y_1, \dots, y_n]^\top$ be the observed targets, and consider a new input x_* with latent value $f_* = f(x_*)$. The joint prior over training and test outputs is

$$\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} K_{\mathbf{xx}} + \sigma^2 I & K_{\mathbf{x}x_*} \\ K_{x_*\mathbf{x}} & K_{x_*x_*} \end{bmatrix}\right),$$

where $K_{\mathbf{x}x_*} = [k(x_1, x_*), \dots, k(x_n, x_*)]^\top$.

Conditioning on the observed data gives the posterior distribution:

$$p(f_* | \mathbf{x}, \mathbf{y}, x_*) = \mathcal{N}(\mu_*, \sigma_*^2),$$

with

$$\mu_* = K_{x_*\mathbf{x}} [K_{\mathbf{xx}} + \sigma^2 I]^{-1} \mathbf{y}, \quad \sigma_*^2 = k(x_*, x_*) - K_{x_*\mathbf{x}} [K_{\mathbf{xx}} + \sigma^2 I]^{-1} K_{\mathbf{x}x_*}.$$

- μ_* is just a linear combination of the observed targets.

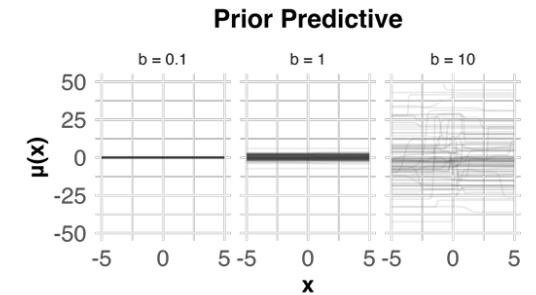
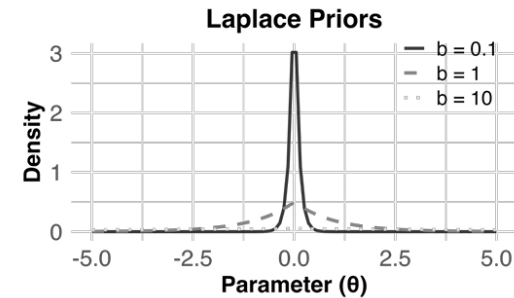
GPs and Deep Learning

- Standard GPs use fixed, pre-defined kernels which operate directly on raw input features. This has the following downsides.
 - Uninformative raw features: For data like images, text or graphs, standard kernels fail to capture the dependencies that underlie the data.
 - High-dimensional inputs: In high-dimensions pairwise distances tend to concentrate, making common kernels nearly constant.
 - Limited adaptability: Once we choose a kernel (and learn its parameters) its inductive bias becomes static. “This lack of flexibility restricts their expressiveness in modern representation learning settings, where adaptive, data-driven feature extraction is essential.”
- **Deep Gaussians Processes** (small number of parameters compared to neural networks): Hierarchical composition of GPs. The output of one GP feeds into another.
- **Deep Kernel Learning:** Use a neural network to learn the kernel.

Note: When a network’s width tends to infinity, its outputs converge to a GP whose kernel is determined by the network architecture (see Neural Tangent Kernel, NIPS 2018, paper with ~5k citations).

Priors for BDL

- In traditional statistics, a prior allows us to represent our initial beliefs before observing any data. It allows us to integrate expert knowledge in our models.
- In deep learning, they usually serve a different purpose. Due to the complexity of the models, it is hard to have any “expert knowledge” on the distribution of the parameters.
- Priors are often selected for practical purposes such as to prevent overfitting.



Informative Priors

- Rare due to model complexity.
- Sometimes applied to specific parameters.
 - Bias term of the output layer.
 - Spike-and-slab priors for the weights associated with covariates in the input layer.

$$p(\theta_{ij}^{(1)} \mid \gamma_i) = \gamma_i \mathcal{N}(\theta_{ij}^{(1)}; \mu_i, \sigma_i^2) + (1 - \gamma_i) \delta_0(\theta_{ij}^{(1)}), \quad \gamma_i \sim \text{Bernoulli}(\rho_i)$$

- Example: If we are predicting heart disease risk, we might assign a high inclusion probability to the weight associated with physical activity and a low one with the weight associated with eye color.

Uninformative Priors?

- Uninformative priors aim to minimize subjective influence. A common choice in classical statistics is a uniform over an interval $[a, b]$.
- **Marginally uniform priors do not always imply an uninformative joint prior especially in high-dimensional BDL models.**
- Reference priors (such as Jeffrey's priors) are designed to maximize the difference between prior and posterior.
- Due to their computational complexity, we can't use them in BDL models!
- In practice, **uninformative priors are rarely used in BDL!**

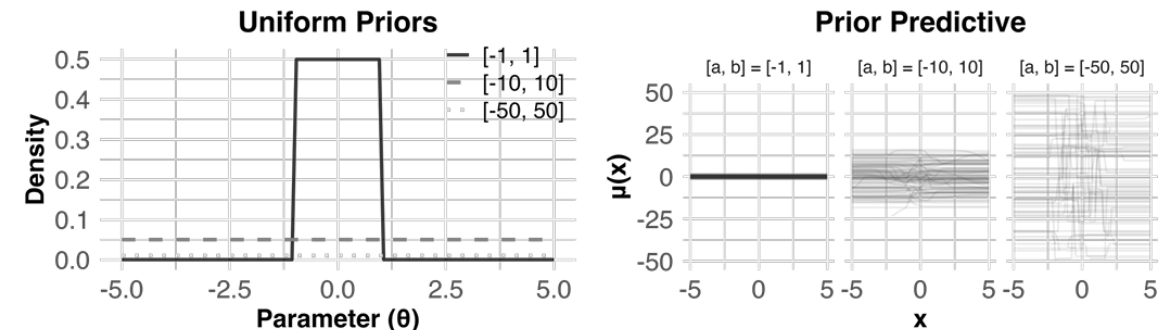


Figure 28.3 Uniform priors for single parameters (left) and the corresponding prior predictive distributions (right), showing a wide range of possible BDL functions with concentration effects for narrow intervals and non-uniform prior predictive.

Regularization Priors

- They constrain model complexity to prevent overfitting.
- The zero-centered isotropic Gaussian prior corresponds to ℓ_2 regularization in MAP estimation.
- The Laplace prior corresponds to ℓ_1 regularization in MAP estimation (encourage parameters to be zero or near-zero).

$$p(\boldsymbol{\theta}) = \prod_{i=1}^P \mathcal{N}(\theta_i; 0, \sigma^2)$$

$$p(\boldsymbol{\theta} \mid \mathcal{M}) = \prod_{i=1}^P \text{Laplace}(\theta_i; 0, b)$$

Priors

- Prior predictive behavior depends on the architecture of the model, particularly the **activation functions**.
- It is advised to conduct prior predictive checks on simplified versions of the model, to avoid prior misspecification.
- Figure shows the prior predictive distribution for Gaussian priors.

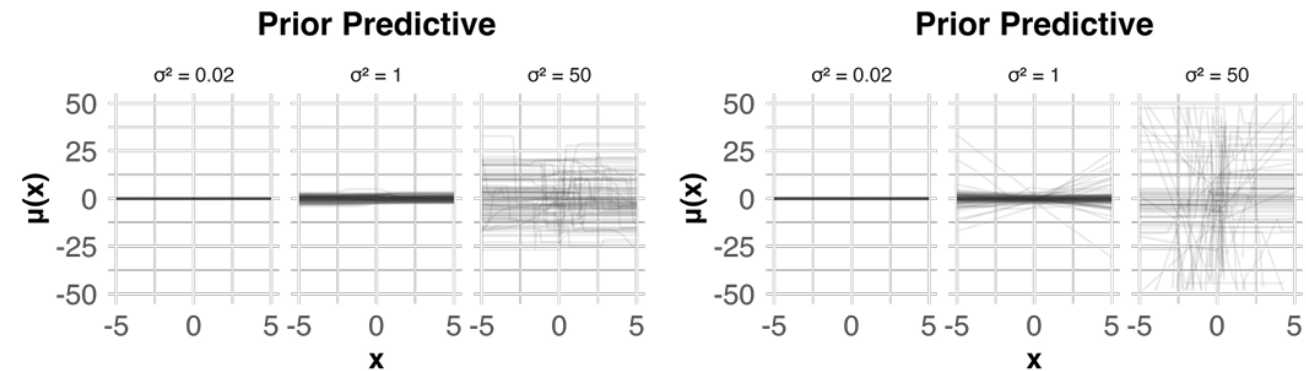
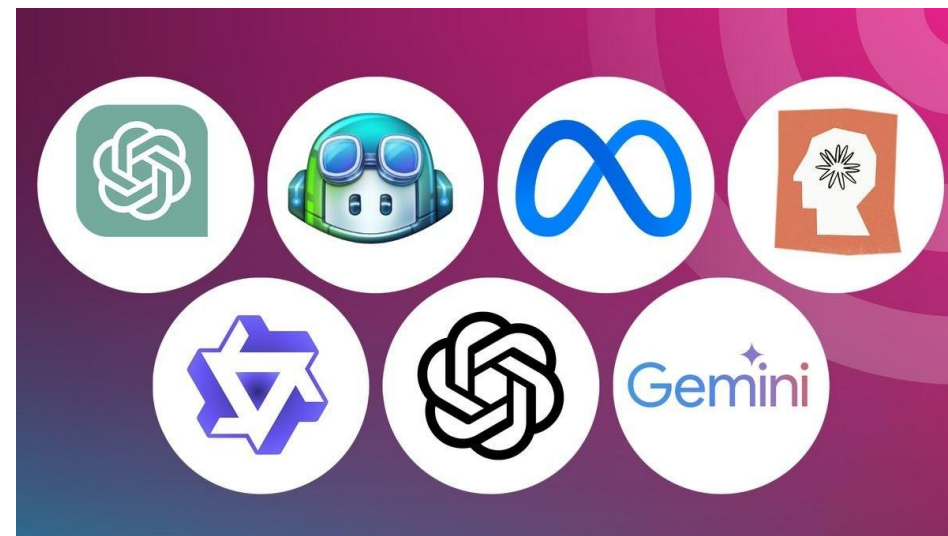


Figure 28.5 Prior predictive distributions for tanh (left) and ReLU (right) activations.

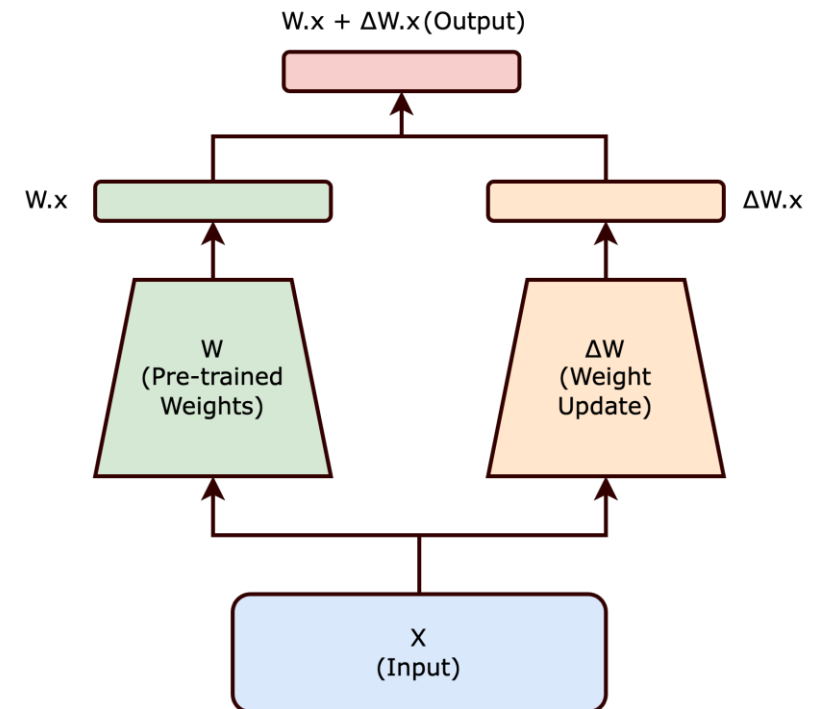
Bayesian LLMs

- LLMs are amazing but have some noticeable limitations.
- The most common of them is known as “hallucinations”, i.e. **confidently** generating incorrect outputs.
- If we had Bayesian LLMs, we would be able to quantify uncertainty. Prediction with model averaging should be more reliable (biases of single models are averaged out).
- Standard Bayesian methods do not work!
 - For training MCMC struggles with slow mixing, VI has high gradient variance due to the dimensionality of the posterior.
 - For inference BMA is very slow.
 - The motivation for being Bayesian is unclear; we have huge training datasets, so the epistemic uncertainty should be low, newer LLMs can defer to external search engines.
- Early works of Bayesian LLMs focused on fine-tuning, where the datasets are small.



Making LoRA Bayesian

- LoRA: Freeze base model weights but add a low rank matrix (which we train).
- $W' = W + \Delta W = W + AB$, $A \in \mathbb{R}^{d \times r}$, $B \in \mathbb{R}^{r \times k}$, $r \ll \min(d, k)$
- LoRA ensemble: Learn multiple ΔW and ensemble the result (either by multiple forward passes or by setting $\Delta W = \sum_{i=1}^N \lambda_i \Delta W_i$).
- Bayesian LoRA: Treat LoRA parameters as random variables and learn a posterior distribution over them.



Improved Variational Online Newton

- Uses VI – places a Normal variational approximation over each parameter (so we have to learn double the parameters).
- The computational costs (for training) are nearly identical to Adam, but its predictive uncertainty is better.
- We want to minimize the following objective:

$$\mathcal{L}(q) = \lambda \mathbb{E}_{q(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] + \mathbb{D}_{\text{KL}}(q(\boldsymbol{\theta}) \| p(\boldsymbol{\theta}))$$

- **Natural gradients** perform gradient descent using the geometry of probability distributions, measuring distances with KL divergence rather than Euclidean distance. They are computed by **preconditioning the ordinary gradient with the inverse Fisher information matrix.**

$$\nabla^{\text{nat}} = F^{-1} \nabla$$

Algorithm 1 Improved Variational Online Newton (IVON). Hyperparameter setting is described in App. [A](#).

Require: Learning rates $\{\alpha_t\}$, weight-decay $\delta > 0$.

Require: Momentum parameters $\beta_1, \beta_2 \in [0, 1)$.

Require: Hessian init $h_0 > 0$.

Init: $\mathbf{m} \leftarrow (\text{NN-weights})$, $\mathbf{h} \leftarrow h_0$, $\mathbf{g} \leftarrow 0$, $\lambda \leftarrow N$.

Init: $\sigma \leftarrow 1/\sqrt{\lambda(\mathbf{h} + \delta)}$.

Optional: $\alpha_t \leftarrow (h_0 + \delta)\alpha_t$ for all t .

```

1: for  $t = 1, 2, \dots$  do
2:    $\hat{\mathbf{g}} \leftarrow \hat{\nabla} \bar{\ell}(\boldsymbol{\theta})$ , where  $\boldsymbol{\theta} \sim q$ 
3:    $\hat{\mathbf{h}} \leftarrow \hat{\mathbf{g}} \cdot (\boldsymbol{\theta} - \mathbf{m}) / \sigma^2$ 
4:    $\mathbf{g} \leftarrow \beta_1 \mathbf{g} + (1 - \beta_1) \hat{\mathbf{g}}$ 
5:    $\mathbf{h} \leftarrow \beta_2 \mathbf{h} + (1 - \beta_2) \hat{\mathbf{h}} + \frac{1}{2} (1 - \beta_2)^2 (\mathbf{h} - \hat{\mathbf{h}})^2 / (\mathbf{h} + \delta)$ 
6:    $\bar{\mathbf{g}} \leftarrow \mathbf{g} / (1 - \beta_1^t)$ 
7:    $\mathbf{m} \leftarrow \mathbf{m} - \alpha_t (\bar{\mathbf{g}} + \delta \mathbf{m}) / (\mathbf{h} + \delta)$ 
8:    $\sigma \leftarrow 1/\sqrt{\lambda(\mathbf{h} + \delta)}$ 
9: end for
10: return  $\mathbf{m}, \sigma$ 
```

Improved Variational Online Newton

- IVON estimates curvature using a reparameterization-based Gauss–Newton approximation, which enables efficient natural-gradient updates!
- IVON@mean uses the mean of the weights.
- IVON uses a model average with 64 samples drawn from the approximation of the posterior.
- Still, lots of disadvantages:
 - Posterior inference requires multiple forward passes.
 - More hyperparameters to tune.

Dataset / Model	Method		Top-1 Acc. \uparrow	Top-5 Acc. \uparrow	NLL \downarrow	ECE \downarrow	Brier \downarrow
ImageNet ResNet-50 (26M params)	AdamW	(2%)	75.16 \pm 0.14	92.37 \pm 0.03	1.018 \pm 0.003	0.066 \pm 0.002	0.349 \pm 0.002
	SGD	(1%)	76.63 \pm 0.45	93.21 \pm 0.25	0.917 \pm 0.026	0.038 \pm 0.009	0.326 \pm 0.006
	IVON@mean		77.30 \pm 0.08	93.58 \pm 0.05	0.884 \pm 0.002	0.035 \pm 0.002	0.316 \pm 0.001
	IVON		77.46 \pm 0.07	93.68 \pm 0.04	0.869 \pm 0.002	0.022 \pm 0.002	0.315 \pm 0.001
TinyImageNet ResNet-18 (11M params)	AdamW	(15%)	47.33 \pm 0.90	71.54 \pm 0.95	6.823 \pm 0.235	0.421 \pm 0.008	0.913 \pm 0.018
	SGD	(1%)	61.39 \pm 0.18	82.30 \pm 0.22	1.811 \pm 0.010	0.138 \pm 0.002	0.536 \pm 0.002
	IVON@mean		62.41 \pm 0.15	83.77 \pm 0.18	1.776 \pm 0.018	0.150 \pm 0.005	0.532 \pm 0.002
	IVON		62.68 \pm 0.16	84.12 \pm 0.24	1.528 \pm 0.010	0.019 \pm 0.004	0.491 \pm 0.001
TinyImageNet PreResNet-110 (4M params)	AdamW	(11%)	50.65 \pm 0.0*	74.94 \pm 0.0*	4.487 \pm 0.0*	0.357 \pm 0.0*	0.812 \pm 0.0*
	SGD	(2%)	59.39 \pm 0.50	81.34 \pm 0.30	2.040 \pm 0.040	0.176 \pm 0.006	0.577 \pm 0.007
	IVON@mean		60.85 \pm 0.39	83.89 \pm 0.14	1.584 \pm 0.009	0.053 \pm 0.002	0.514 \pm 0.003
	IVON		61.25 \pm 0.48	84.13 \pm 0.17	1.550 \pm 0.009	0.049 \pm 0.002	0.511 \pm 0.003
CIFAR-100 ResNet-18 (11M params)	AdamW	(11%)	64.12 \pm 0.43	86.85 \pm 0.51	3.357 \pm 0.071	0.278 \pm 0.005	0.615 \pm 0.008
	SGD	(1%)	74.46 \pm 0.17	92.66 \pm 0.06	1.083 \pm 0.007	0.113 \pm 0.001	0.376 \pm 0.001
	IVON@mean		74.51 \pm 0.24	92.74 \pm 0.19	1.284 \pm 0.013	0.152 \pm 0.003	0.399 \pm 0.002
	IVON		75.14 \pm 0.34	93.30 \pm 0.19	0.912 \pm 0.009	0.021 \pm 0.003	0.344 \pm 0.003
CIFAR-100 PreResNet-110 (4M params)	AdamW	(10%)	65.88 \pm 0.84	88.34 \pm 0.56	2.893 \pm 0.088	0.258 \pm 0.006	0.578 \pm 0.014
	SGD	(2%)	74.19 \pm 0.11	92.41 \pm 0.14	1.204 \pm 0.012	0.137 \pm 0.002	0.393 \pm 0.004
	IVON@mean		75.23 \pm 0.23	93.45 \pm 0.16	1.149 \pm 0.010	0.136 \pm 0.002	0.380 \pm 0.003
	IVON		75.81 \pm 0.18	93.93 \pm 0.19	0.884 \pm 0.007	0.030 \pm 0.003	0.336 \pm 0.001

Table 1: IVON improves both accuracy and uncertainty over SGD and AdamW. Improvements in accuracy by IVON are shown in red. The performance of AdamW is not good on the smaller datasets likely due to overfitting when training for 200 epochs. IVON does not have this issue. Additional results are in the appendix in Tables 9 to 11 and Table 13.

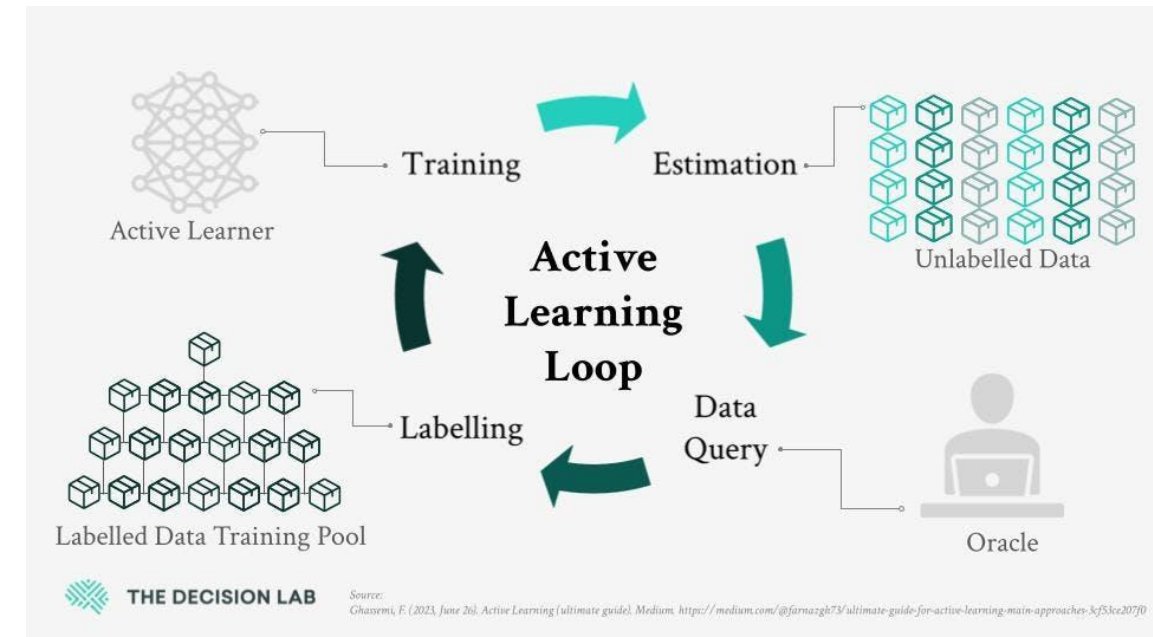
BDL in Practice

- In practice when would you use something like IVON?
- Medical image diagnosis (image classification).
- For each parameter we sample perform a forward pass.
- Build a credible interval to account for uncertainty!

DO NOT USE BDL IF YOU DO NOT CARE ABOUT UNCERTAINTY.

Active Learning

- **Which data do we want to acquire?**
- **Key Idea:** Labeled data are expensive to obtain. How do we choose which data points we want labeled?
- Toy active learning example:
 - Train a small model on a **labeled** dataset.
 - Pass a large pool of **unlabeled** data through the trained model.
 - Choose the most informative unlabeled data and label them (e.g., output probabilities close to 0.5).
 - Retrain model.



Thank you for
listening!
Questions?

