# Developments in Inference with Linearised Neural Networks

Riccardo Barbano, Javier Antorán

UNIVERSITY OF
CAMBRIDGE

# Reading group outline

1. **Preliminaries**: probabilistic inference in neural networks and the linearised Laplace method

2. **Paper overview**: "Improving predictions of Bayesian neural networks via local linearization"

3. **Informal discussion**

# Preliminaries

## Overconfidence

Training on CIFAR10 – Test on SVHN



Dog (100%)    Bird (100%)    Airplane (100%)

https://vitalab.github.io/article/2019/07/11/overconfident.html

## Model Selection

1 Hidden Layer          5 Hidden Layer          20 Hidden Layer

1. Place a prior distribution $\pi(\theta)$ over NN parameters.
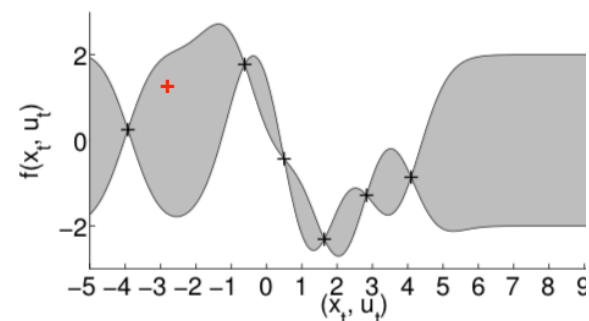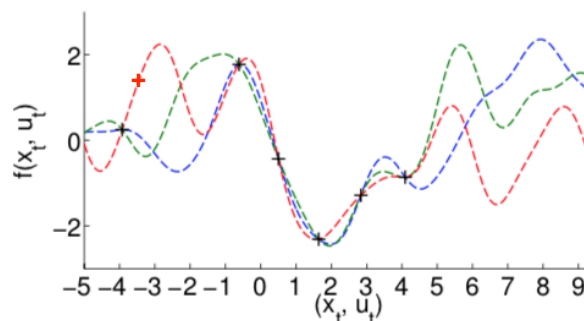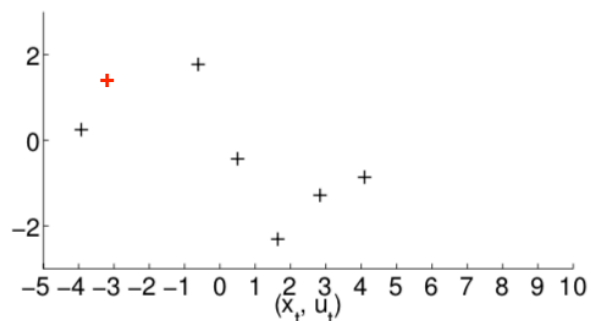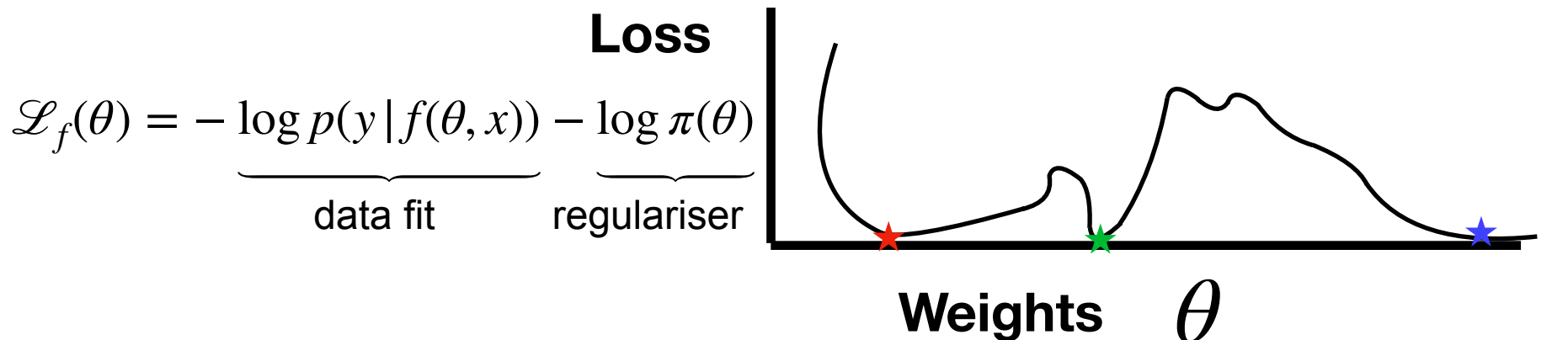
2. Define some likelihood function $p(y\,|\,f(\theta, x))$ to characterise the agreement of the NN function $f(\theta, \cdot)$ with the observations $(y, x)$

3. Update the weight distribution using Bayes' rule



$$\tilde{\theta} \in \operatorname{argmax}_\theta \log p(y\,|\,f(\theta, x)) + \log \pi(\theta)$$

$$\underbrace{\qquad\qquad}_{\text{data fit}} \quad \underbrace{\qquad}_{\text{regulariser}}$$

$$p(\theta\,|\,x, y) = \frac{p(y\,|\,f(\theta, x))\pi(\theta)}{p(y\,|\,x)}$$

# Preliminaries: uncertainty estimation

$$\mathscr{L}_f(\theta) = -\underbrace{\log p(y \mid f(\theta, x))}_{\text{data fit}} - \underbrace{\log \pi(\theta)}_{\text{regulariser}}$$
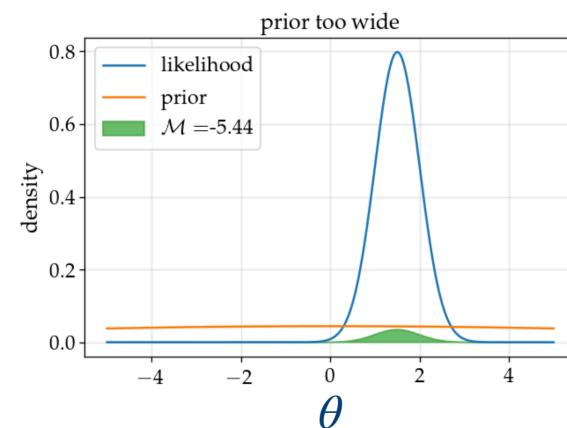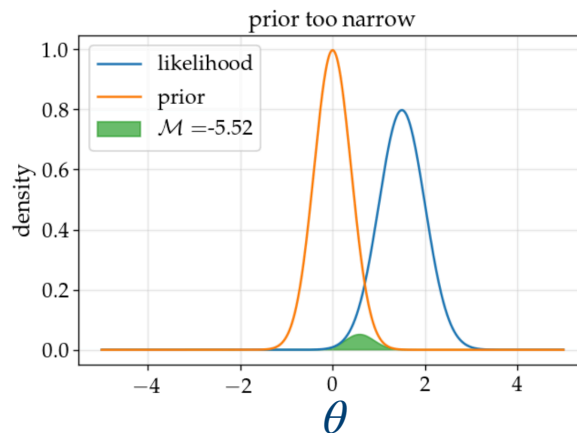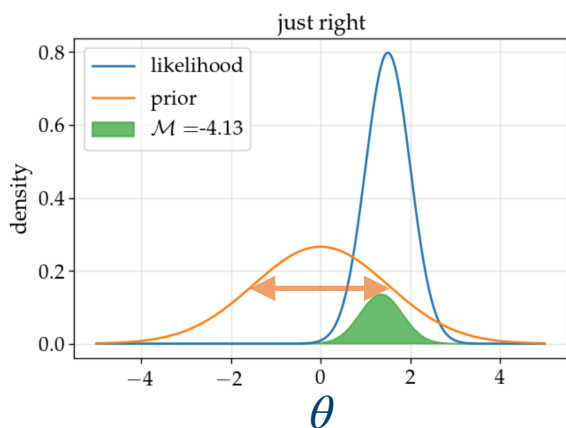
**Loss**



**Weights** $\theta$



$$p(\theta \mid y, x) = \frac{1}{\exp(\mathscr{M})} \exp(-\mathscr{L}_f(\theta)) \qquad f(\theta, \cdot), \ \theta \sim p(\theta \mid y, x)$$
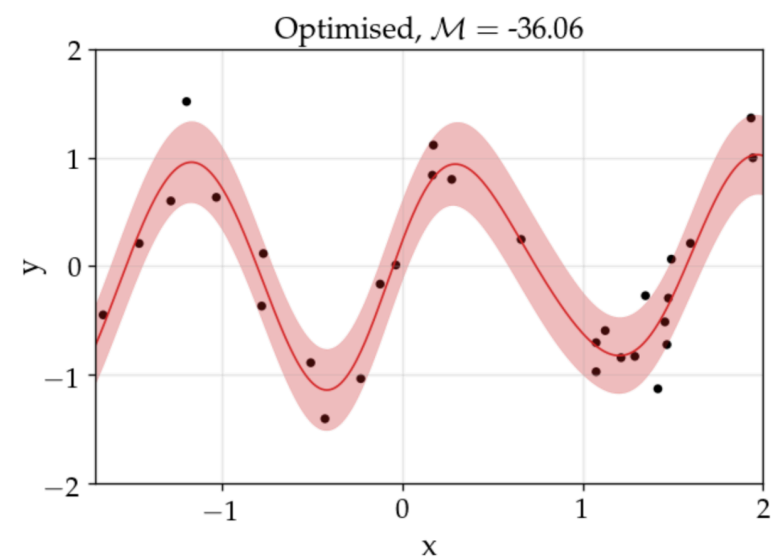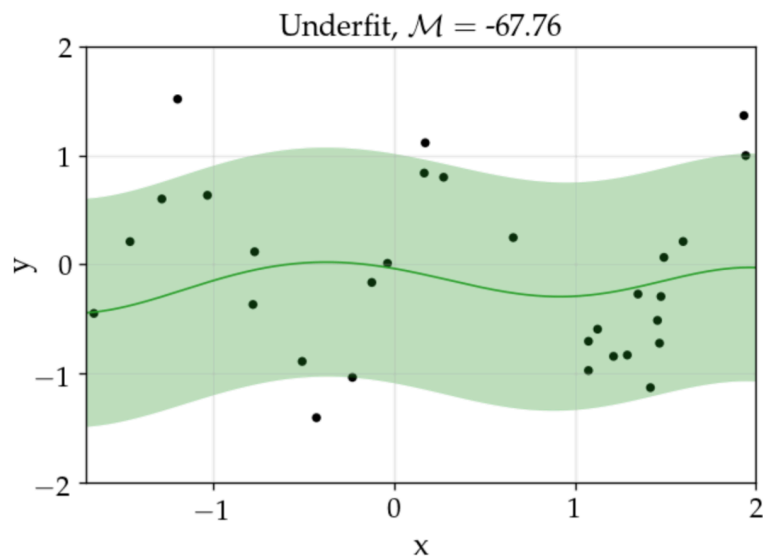
$$\mathscr{L}_f(\theta) = -\log p(y \,|\, f(\theta, x)) - \log \pi(\theta) \quad p(\theta \,|\, y, x) = \frac{1}{\exp(\mathscr{M})} \exp(-\mathscr{L}_f(\theta))$$

The normalisation constant, $\mathscr{M}$, is the **marginal likelihood**, or **model evidence**. It is the probability that our observations where generated by our prior. It provides an objective for hyperparameter selection without the need for validation data.

$$\mathscr{M} = \log p(y \,|\, x) \;\; = \log \int p(y \,|\, f(\theta, x)) d\pi \;\; = \log \int \exp(-\mathscr{L}_f(\theta)) d\nu$$

# Preliminaries: automatic Occam's razor



All possible data sets

# Preliminaries: the Laplace approximation
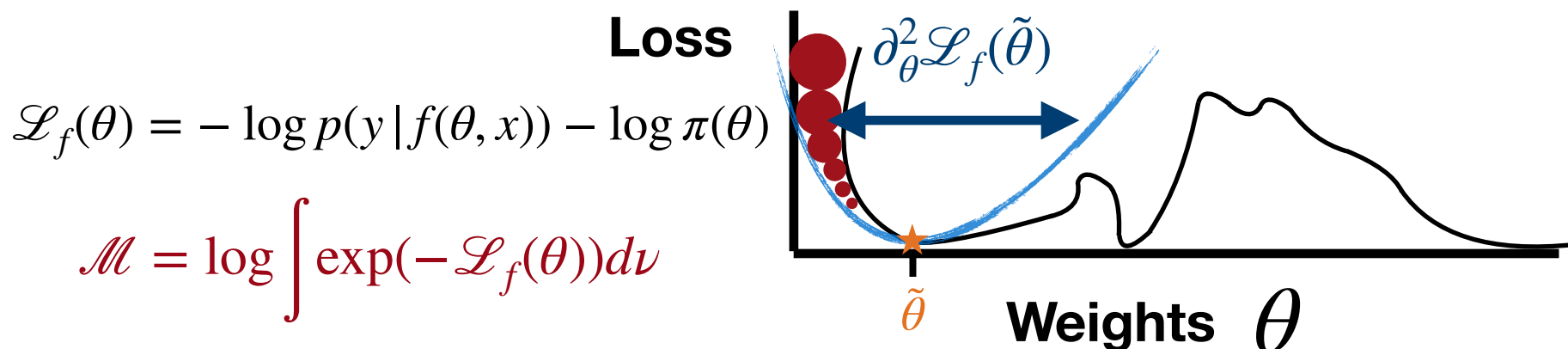
$$\mathcal{L}_f(\theta) = -\log p(y\,|\,f(\theta, x)) - \log \pi(\theta)$$

$$\mathcal{M} = \log \int \exp(-\mathcal{L}_f(\theta))d\nu$$

For NNs this integral is intractable

**Loss** $\partial_\theta^2 \mathcal{L}_f(\tilde{\theta})$

$\tilde{\theta}$ **Weights** $\theta$

**Idea**: Find a mode of $\mathcal{L}_f$: $\tilde{\theta}$ and perform 2-order Taylor expansion

$$\mathcal{G}_f(\theta) = \mathcal{L}_f(\tilde{\theta}) + ||\theta - \tilde{\theta}||^2_{\partial_\theta^2 \mathcal{L}_f(\tilde{\theta})}$$

By inspection, $\exp(-\mathcal{G}_{f,\tilde{\theta}}(\theta))$ is proportional to $\mathcal{N}(\tilde{\theta},\ (\partial_\theta^2 \mathcal{L}_f(\tilde{\theta}))^{-1})$ where

$$\partial_\theta^2 \mathcal{L}_f(\tilde{\theta})) = \partial_\theta^2 \log p(y\,|\,f(\tilde{\theta}, x)) + \partial_\theta^2 \log \pi(\tilde{\theta}) \qquad \pi(\theta) \to \mathcal{N}(\theta;\ 0, \Lambda^{-1})$$

**Issue**: A lot of mass falls in low density region, leading to bad predictions

PAPER DISCUSSION:

IMPROVING PREDICTIONS OF BAYESIAN NEURAL NETWORKS VIA LOCAL LINEARIZATION