

# PROBABILISTIC MACHINE LEARNING

## LECTURE 19

### USES FOR UNCERTAINTY IN DEEP LEARNING II

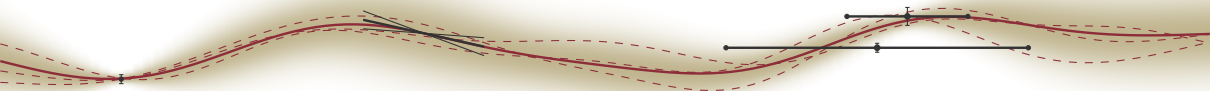
Philipp Hennig

6 July 2023

EBERHARD KARLS  
UNIVERSITÄT  
TÜBINGEN



FACULTY OF SCIENCE  
DEPARTMENT OF COMPUTER SCIENCE  
CHAIR FOR THE METHODS OF MACHINE LEARNING



# Recap: Deep networks *are* GPs

Linearized networks and Laplace approximations: From deep learning to GPs, in four easy steps

1. Realise that the loss is a **negative log-posterior**

$$\mathcal{L}(\boldsymbol{\theta}) = \left( \underbrace{\frac{1}{N} \sum_{i=1}^N \ell(y_i; \overbrace{f(x_i, \boldsymbol{\theta})}^{\text{deep net}})}_{\text{empirical risk}} + \underbrace{r(\boldsymbol{\theta})}_{\text{regularizer}} \right) = - \sum_{i=1}^N \log p(y | \boldsymbol{\theta}) - \log p(\boldsymbol{\theta}) = - \log p(\boldsymbol{\theta} | y) + \text{const.}$$

2. Train the deep net *as usual* to find  $\boldsymbol{\theta}_* = \arg \max_{\boldsymbol{\theta} \in \mathbb{R}^D} p(\boldsymbol{\theta} | y)$
3. At  $\boldsymbol{\theta}_*$ , compute a **Laplace approximation** of the log-posterior, with  $\Psi := \nabla \nabla^\top \log p(\boldsymbol{\theta}_* | y)$

$$\log p(\boldsymbol{\theta} | y) + \text{const.} = \mathcal{L}(\boldsymbol{\theta}) \approx \mathcal{L}(\boldsymbol{\theta}_*) + 1/2 (\boldsymbol{\theta} - \boldsymbol{\theta}_*)^\top \Psi (\boldsymbol{\theta} - \boldsymbol{\theta}_*) = \log \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\theta}_*, -\Psi^{-1})$$

4. Linearize  $f(x, \boldsymbol{\theta})$  around  $\boldsymbol{\theta}_*$ , with  $[J(x)]_{ij} = \frac{\partial f_i(x, \boldsymbol{\theta}_*)}{\partial \theta_j}$  as  $f(x, \boldsymbol{\theta}) \approx f(x, \boldsymbol{\theta}_*) + J(x, \boldsymbol{\theta}_*)(\boldsymbol{\theta} - \boldsymbol{\theta}_*)$

$$\text{thus } p(f(\bullet) | \mathcal{D}) = \int p(f | w) dp(w) \approx \mathcal{GP}(f(\bullet); f(\bullet, \boldsymbol{\theta}_*), -J(\bullet)\Psi^{-1}J(\circ)) \quad \text{with}$$

$$\mathbb{E}(f(\bullet)) = f(\bullet, \boldsymbol{\theta}_*) \quad \text{the trained net as the mean function}$$

$$\text{cov}(f(\bullet), f(\circ)) = -J(\bullet)\Psi^{-1}J(\circ)^\top \quad \text{the Laplace tangent kernel as the covariance function}$$

Computing the exact Hessian  $\Psi$  is  $\mathcal{O}(ND^2)$ , and inverting it is  $\mathcal{O}(D^3)$ .

Ideas for Approximations:

- ▶ Sub-sample the dataset ( $\mathcal{O}(MD^2)$  with  $M \ll N$ )
- ▶ structural approximations to the Hessian:
  - ▶ diagonal approximation:  $\mathcal{O}(D)$  (inverse  $\mathcal{O}(D)$ )
  - ▶ last-layer approximation:  $\mathcal{O}(D_L^2)$  (inverse  $\mathcal{O}(D_L^3)$ )
  - ▶ Kronecker factorized approximate curvature (KFAC):  $\Psi \approx \text{diag}([\Lambda_\ell \otimes \Omega_\ell]_{\ell=1,\dots,L})$   
with  $\Lambda_\ell \in \mathbb{R}^{\text{in}_\ell \times \text{in}_\ell}$ ,  $\Omega_\ell \in \mathbb{R}^{\text{out}_{\ell-1} \times \text{out}_{\ell-1}}$  and thus inverse  $\mathcal{O}\left(\sum_\ell \text{in}_\ell^3 + \text{out}_{\ell-1}^3\right)$
  - ▶ Generalized Gauss-Newton (homework this week):  $\Psi \approx \alpha I + GG^\top$  with  $G \in \mathbb{R}^{D \times M}$
  - ▶ approximate eigenvalue decompositions using the Lanczos algorithm (cf. Lecture 13)

# What's not to like about deep learning?

Reasons to look to probability theory for answers



Review from Lecture 16

What people **don't** like about deep learning:

- ▶ deep learning has certain conceptual pathologies, leading to adversarial and out-of-distribution brittleness
- ▶ Training a deep net is fiddly, and requires *many* choices
- ▶ Once the model is trained, it's unclear how to *update* it if new data arrives

Geometric / probabilistic interpretation of deep learning can help address these issues.



$$\theta_* = \arg \min_{\theta \in \mathbb{R}^D} \underbrace{\mathcal{L}(\theta)}_{\text{Loss}} = \arg \min_{\theta \in \mathbb{R}^D} \left( \underbrace{\frac{1}{N} \sum_{i=1}^N \ell(y_i; \overbrace{f(x_i, \theta)}^{\text{deep net}})}_{\text{empirical risk}} + \underbrace{r(\theta)}_{\text{regularizer}} \right)$$

**Theorem:** [Hein et al. 2019] If  $f(x, \theta)$  is a ReLU network and  $\ell = \sigma$  is the cross-entropy then, for almost any  $x \in \mathbb{R}^d$  and  $\varepsilon > 0$ , there exists an  $\alpha > 0$  and a class  $k \in \{1, \dots, K\}$  such that for  $z = \alpha x$

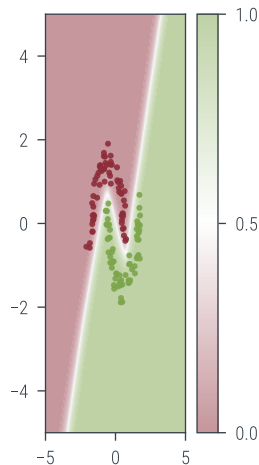
$$\frac{e^{f_k(x)}}{\sum_{r=1}^K e^{f_r(z)}} \geq 1 - \varepsilon \quad \text{and} \quad \lim_{\alpha \rightarrow \infty} \frac{e^{f_k(\alpha x)}}{\sum_{r=1}^K e^{f_r(\alpha x)}} = 1$$

*Far from the data, the model is always confident.*

Intuition: ReLU networks are piecewise linear, so far from the data,

$f_i(x, \theta) = A_i(\theta)x$  and  $f_{\arg \max_i A_i}(x) - f_j = \mathcal{O}(\|x\|)$ , so

$\lim_{\|x\| \rightarrow \infty} \text{softmax}(f(x)) = 1$ .



Sketch: (Proof in Kristiadi, Hein, Hennig, ICML 2020)

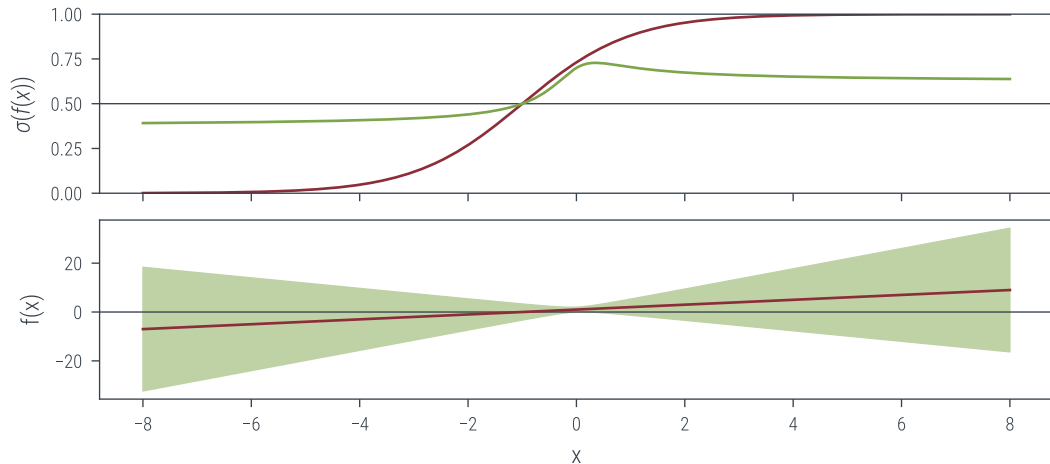
- Far from the data,  $f_i(\mathbf{x}, \boldsymbol{\theta}) = A_i(\boldsymbol{\theta})\mathbf{x}$
- Consider *any* Gaussian measure on the weights  $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}, \boldsymbol{\theta}_*, \Psi)$  (e.g. from Laplace approximation). The associated push-forward GP is

$$p(f(\mathbf{x})) = \mathcal{GP}(f(\mathbf{x}), f(\mathbf{x}, \boldsymbol{\theta}_*), J(\mathbf{x}, \boldsymbol{\theta}_*)\Psi J(\mathbf{x}, \boldsymbol{\theta}_*))$$
$$\|\mathbf{x}\| \xrightarrow{\infty} \mathcal{GP}(f(\mathbf{x}); A(\boldsymbol{\theta}_*)\mathbf{x}, \mathbf{x}^\top A'(\boldsymbol{\theta}_*)\Psi A'(\boldsymbol{\theta}_*)^\top \mathbf{x})$$

- The class prediction of this network is

$$\mathbb{E}(\sigma(f(\mathbf{x}))) = \int \sigma(f(\mathbf{x}))p(f(\mathbf{x})) d\mathbf{x} \approx \sigma\left(\frac{\mathbb{E}(f(\mathbf{x}))}{\sqrt{1 + \frac{\pi}{8} \text{var}(f(\mathbf{x}))}}\right) \|\mathbf{x}\| \xrightarrow{\infty} \sigma\left(\frac{\mathcal{O}(\|\mathbf{x}\|)}{\sqrt{1 + \mathcal{O}(\|\mathbf{x}\|^2)}}\right)$$
$$\rightarrow \sigma(c) < 1 \quad \text{and} \quad > 0$$

# Being Bayesian, even just a bit, fixes overconfidence



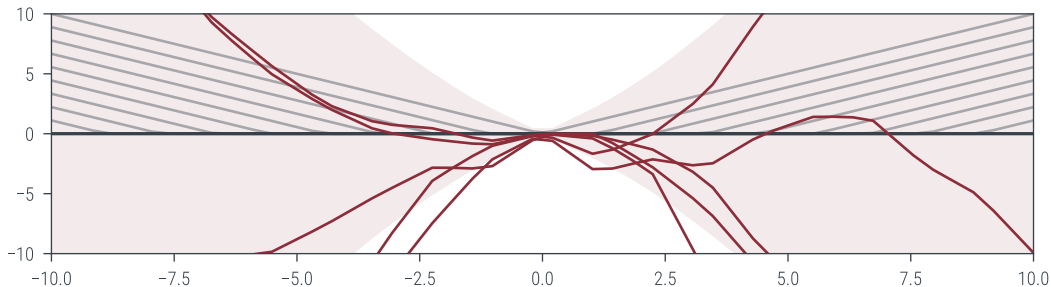
# Finite networks have finite capacity

a role for nonparametric modeling in deep learning



Kristiadi, Hein, Hennig, NeurIPS 2021

- ▶ To achieve  $\mathbb{E}(\sigma(f(\mathbf{x}))) = 1/C$  (“calibrated confidence”) at  $\|\mathbf{x}\| \rightarrow \infty$ , we need  $\text{var } f(\mathbf{x}) = \omega(\|\mathbf{x}\|^2)$ .
- ▶ This can *not* be achieved with finitely many ReLU features in the network (there will always be a “last ReLU to switch on”, and  $f(\mathbf{x}) = A\mathbf{x}$  beyond that).
- ▶ But remember from Lecture 9 there’s a limit process  $\mathcal{GP}(0, k_{IWP})$  with  $k_{IWP}(x, x) = \mathcal{O}(x^3)$





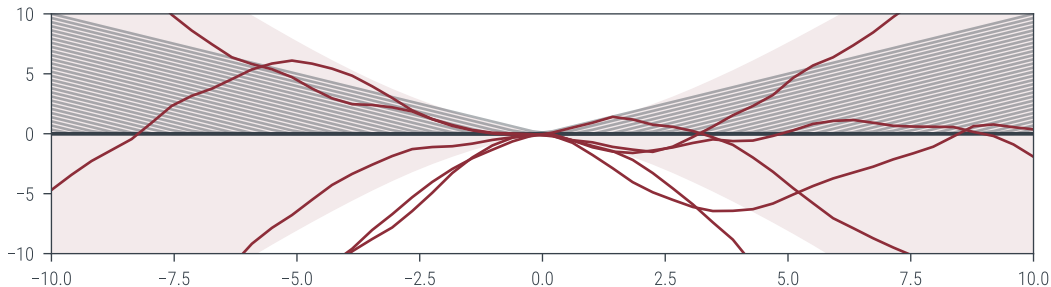
# Finite networks have finite capacity

a role for nonparametric modeling in deep learning



Kristiadi, Hein, Hennig, NeurIPS 2021

- ▶ To achieve  $\mathbb{E}(\sigma(f(\mathbf{x}))) = 1/C$  ("calibrated confidence") at  $\|\mathbf{x}\| \rightarrow \infty$ , we need  $\text{var } f(\mathbf{x}) = \omega(\|\mathbf{x}\|^2)$ .
- ▶ This can *not* be achieved with finitely many ReLU features in the network (there will always be a "last ReLU to switch on", and  $f(\mathbf{x}) = A\mathbf{x}$  beyond that).
- ▶ But remember from Lecture 9 there's a limit process  $\mathcal{GP}(0, k_{IWP})$  with  $k_{IWP}(x, x) = \mathcal{O}(x^3)$



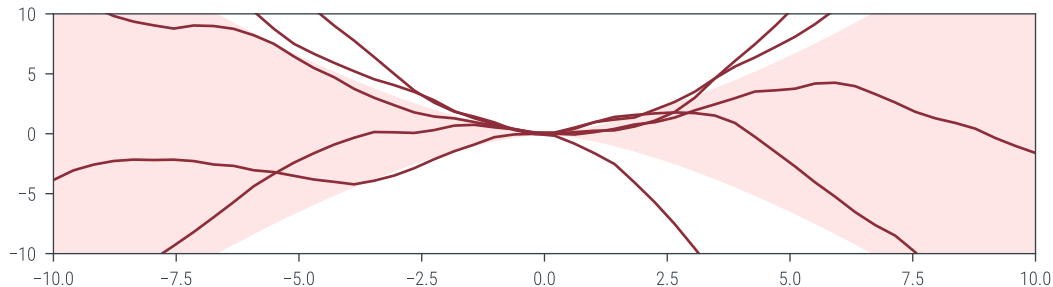


# Finite networks have finite capacity

a role for nonparametric modeling in deep learning

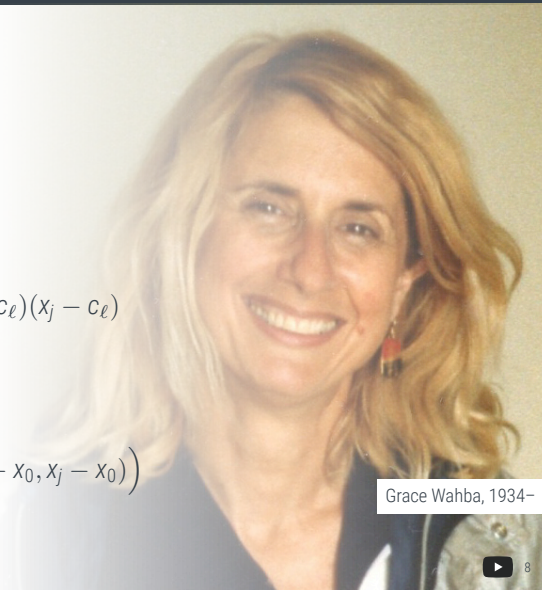
Kristiadi, Hein, Hennig, NeurIPS 2021

- ▶ To achieve  $\mathbb{E}(\sigma(f(\mathbf{x}))) = 1/C$  ("calibrated confidence") at  $\|\mathbf{x}\| \rightarrow \infty$ , we need  $\text{var } f(\mathbf{x}) = \omega(\|\mathbf{x}\|^2)$ .
- ▶ This can *not* be achieved with finitely many ReLU features in the network (there will always be a "last ReLU to switch on", and  $f(\mathbf{x}) = A\mathbf{x}$  beyond that).
- ▶ But remember from Lecture 9 there's a limit process  $\mathcal{GP}(0, k_{IWP})$  with  $k_{IWP}(x, x) = \mathcal{O}(x^3)$





$$\begin{aligned}\phi_i(x) &= \int_{-\infty}^x \theta(\tilde{x} - c_i) d\tilde{x} = \max(0, x - c_i) \\ k(x_i, x_j) &= \frac{\sigma^2(c_{\max} - c_0)}{F} \sum_{\ell} \phi_{\ell}(x_i) \phi_{\ell}(x_j) \\ &= \frac{\sigma^2(c_{\max} - c_0)}{F} \sum_{\ell} \max(0, \min(x_i, x_j) - c_{\ell})(x_i - c_{\ell})(x_j - c_{\ell}) \\ &\rightarrow \sigma^2 \int_{c_0}^{\min(x_i, x_j)} (x_i - c)(x_j - c) dc \\ &= \sigma^2 \left( \frac{1}{3} \min^3(x_i - x_0, x_j - x_0) + \frac{1}{2} |x_i - x_j| \min^2(x_i - x_0, x_j - x_0) \right)\end{aligned}$$



Grace Wahba, 1934–

- Consider a dataset  $X, y$  w.l.o.g. scaled such that,  $\|x\| \leq 1$  for all  $x \in X$ .
- Model classifier function as  $\tilde{f}(x) = f(x) + \hat{f}(x)$  with a ReLU network  $f(x)$  and  $\hat{f}(x) \sim \mathcal{GP}(0, k_{IWP})$ .
- Then (Lemma in op.cit.), for  $0 < \delta < 1$ , if  $\|a\|^2, \|b\|^2 \leq \delta$ , then  $k_{IWP}(a, b) \in \mathcal{O}(\delta^3)$ .

**Proposition:** Assume the ReLU networks weights are distributed as  $\mathcal{N}(\theta, \theta_*, \Sigma)$  and linearize  $f(x) = f(x, \theta_*) + J(x, \theta_*)(\theta - \theta_*)$ . Then the predictive GP posterior on  $\tilde{f}$  has mean and covariance

$$\begin{aligned}\mathbb{E}(\tilde{f}(\bullet) \mid X, y) &= f(\bullet, \theta_*) + k_{\bullet, X} C^{-1} (y - f(X, \theta_*)) \\ \text{var}(\tilde{f}(\bullet) \mid X, y) &= J(\bullet, \theta_*) \Sigma J^\top(\bullet, \theta_*) + k_{\bullet, \bullet} - k_{\bullet, X} C^{-1} k_{X, \bullet}\end{aligned}$$

with  $C := k_{IWP}(X, X) + \sigma^2 I + J(X, \theta_*) \Sigma J^\top(X, \theta_*)$ . If the network is well-trained ( $f(X, \theta_*) \approx y$ ), then **the residual is negligible**. We also have (simplified argument)

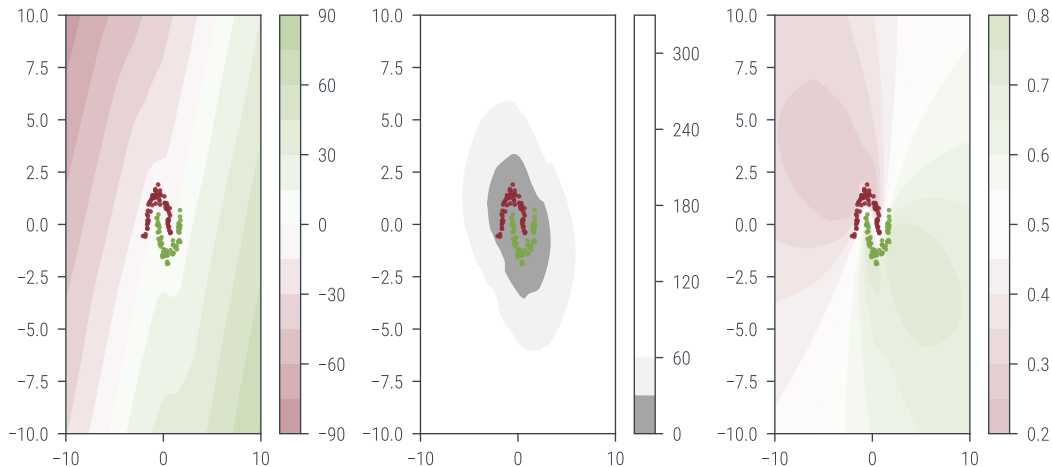
$$k_{\bullet, X} C^{-1} k_{X, \bullet} \leq \lambda_{\max}(C^{-1}) \|k_{\bullet, X}\|^2 = \lambda_{\max}(C^{-1}) \sum_{i=1}^N k_{\bullet, x_i}^2 \leq \lambda_{\max}(C^{-1}) \sum_{i=1}^N k_{\bullet, \bullet} \underbrace{k_{x_i, x_i}}_{\approx 0}$$

Thus:  $p(\tilde{f}(\bullet)) \approx \mathcal{GP}(\tilde{f}(\bullet); f(\bullet, \theta_*), J(\bullet, \theta_*) \Sigma J^\top(\bullet, \theta_*) + k_{\bullet, \bullet})$ .



Adding “infinitely many” untrained units to a deep net

Kristiadi, Hein, Hennig, NeurIPS 2021





Consider a sequence of datasets  $\mathcal{D}_i, i = 1, \dots, T$ . At time  $i$ , previous datasets are no longer available. What do you do?

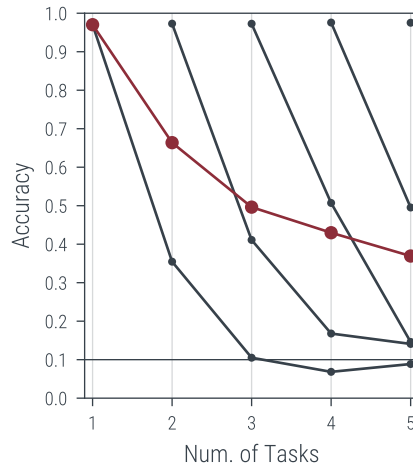
## 1. Just keep training™

(common alternative: "replay" old data

(e.g. <https://openreview.net/forum?id=Q4aAITDgdP>). Which is expensive)

### Setup:

- ▶ Permuted copies of MNIST
- ▶ (784, 100, 10)-ReLU net
- ▶ Adam(lr=10<sup>-3</sup>), 10 epochs each



# Probabilistic inference naturally deals with continual learning

"yesterday's posterior is today's prior"

$$\begin{aligned}
 p(\boldsymbol{\theta} \mid \mathbf{y}_1, y_2) &= \frac{p(\boldsymbol{\theta} \mid \mathbf{y}_1) \cdot p(y_2 \mid \boldsymbol{\theta})}{p(y_2)} && \text{(assuming } y_2 \perp\!\!\!\perp \mathbf{y}_1 \mid \boldsymbol{\theta} \text{)} \\
 -\log p(\boldsymbol{\theta} \mid \mathbf{y}_1, y_2) &= -\log p(y_2 \mid \boldsymbol{\theta}) - \log p(\boldsymbol{\theta} \mid \mathbf{y}_1) \\
 &= \sum_{i=1}^{N_2} \ell(y_{2,i}, f(x_{2,i}, \boldsymbol{\theta})) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_1)^\top \Psi_1(\boldsymbol{\theta} - \boldsymbol{\theta}_1)
 \end{aligned}$$

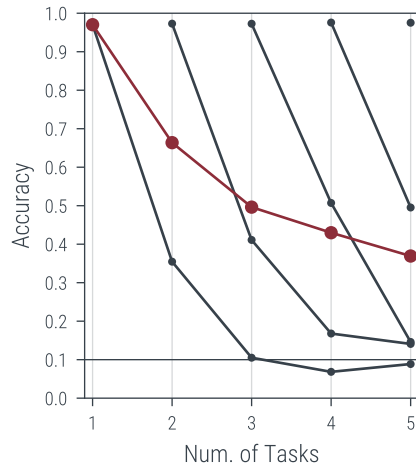


Consider a sequence of datasets  $\mathcal{D}_i, i = 1, \dots, T$ . At time  $i$ , previous datasets are no longer available. What do you do?

1. Just keep training™

## Setup:

- ▶ Permuted copies of MNIST
- ▶ (784, 100, 10)-ReLU net
- ▶ Adam(lr=10<sup>-3</sup>), 10 epochs each





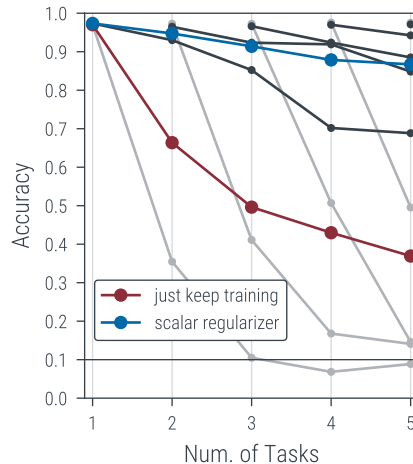


Consider a sequence of datasets  $\mathcal{D}_i, i = 1, \dots, T$ . At time  $i$ , previous datasets are no longer available. What do you do?

1. Just keep training<sup>TM</sup>
2. regularizer  $r(\boldsymbol{\theta}) = \frac{\lambda}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}_{i-1}\|^2$

## Setup:

- ▶ Permuted copies of MNIST
- ▶ (784, 100, 10)-ReLU net
- ▶ Adam(lr=10<sup>-3</sup>), 10 epochs each
- ▶  $\lambda = 5 \cdot 10^{-4}$



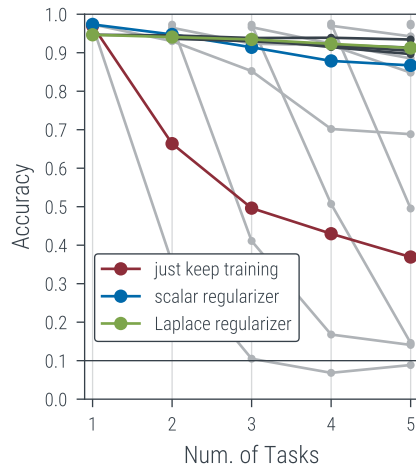


Consider a sequence of datasets  $\mathcal{D}_i, i = 1, \dots, T$ . At time  $i$ , previous datasets are no longer available. What do you do?

1. Just keep training<sup>TM</sup>
2. regularizer  $r(\boldsymbol{\theta}) = \frac{\lambda}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}_{i-1}\|^2$
3. Laplace posterior  $r(\boldsymbol{\theta}) = \frac{\lambda}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_{i-1})^\top \Psi (\boldsymbol{\theta} - \boldsymbol{\theta}_{i-1})$

## Setup:

- ▶ Permuted copies of MNIST
- ▶ (784, 100, 10)-ReLU net
- ▶ Adam(lr=10<sup>-3</sup>), 10 epochs each
- ▶  $\lambda = 5 \cdot 10^{-4}$



## Uncertainty in Deep Learning

- ▶ fixes (asymptotic and local) overconfidence
- ▶ yields the functionality for continual learning
- ▶ many other applications not discussed here

Laplace approximations turn deep networks into GPs, inheriting all functionality of GPs

Please cite this course, as

```
@techreport{Tuebingen_ProbML23,
  title =
    {Probabilistic Machine Learning},
  author = {Hennig, Philipp},
  series = {Lecture Notes
            in Machine Learning},
  year = {2023},
  institution = {Tübingen AI Center}}
```

