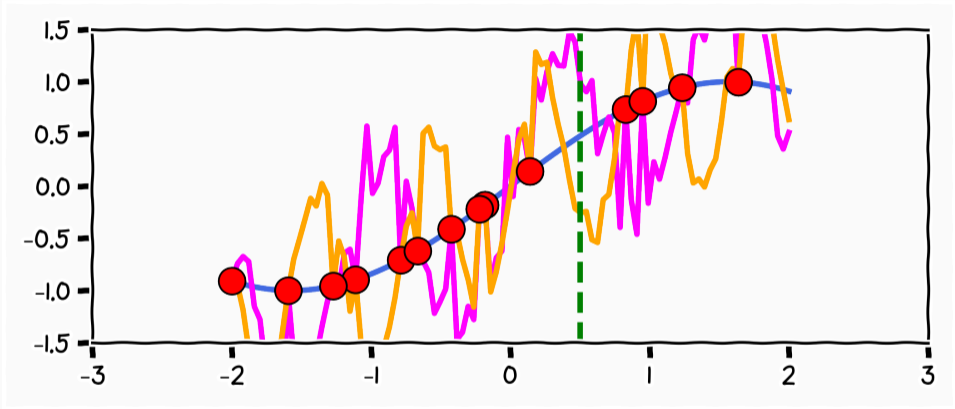


Approximate Bayesian Inference of Composite Functions

Carl Henrik Ek

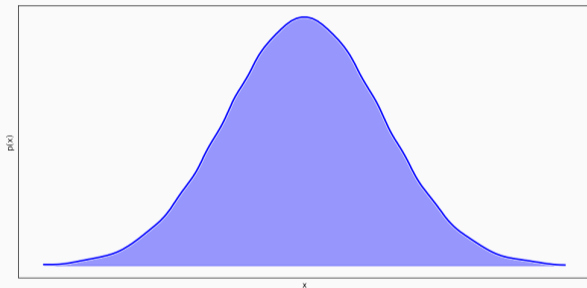
September 9, 2025

<http://carlhenrik.com>

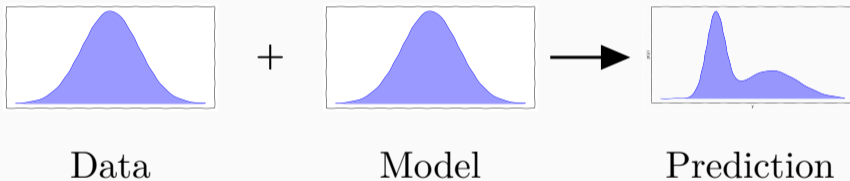


What is Machine Learning

data + model \rightarrow prediction

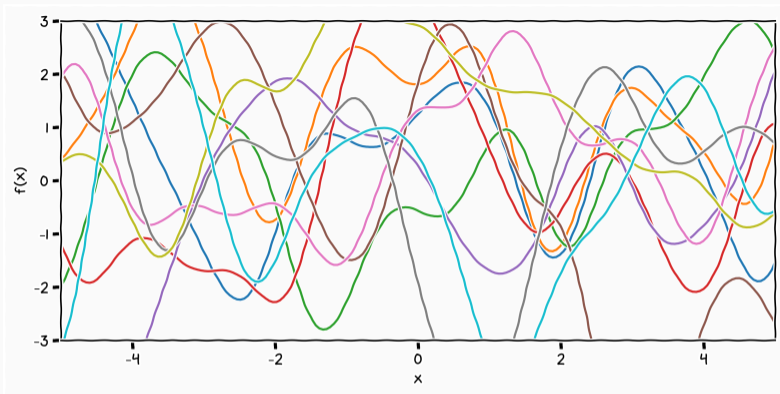


Knowledge (Uncertainty) Propagation





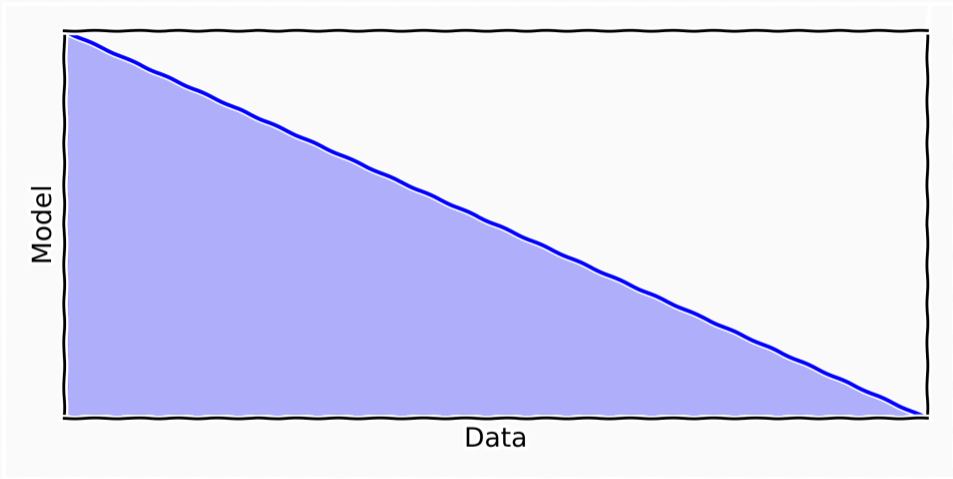
$$p(\mathcal{D}) = \int p(\mathcal{D} \mid f)p(f)\mathrm{d}f$$



$$p(f) = \mathcal{GP}(\mu(x), k(x, x'))$$

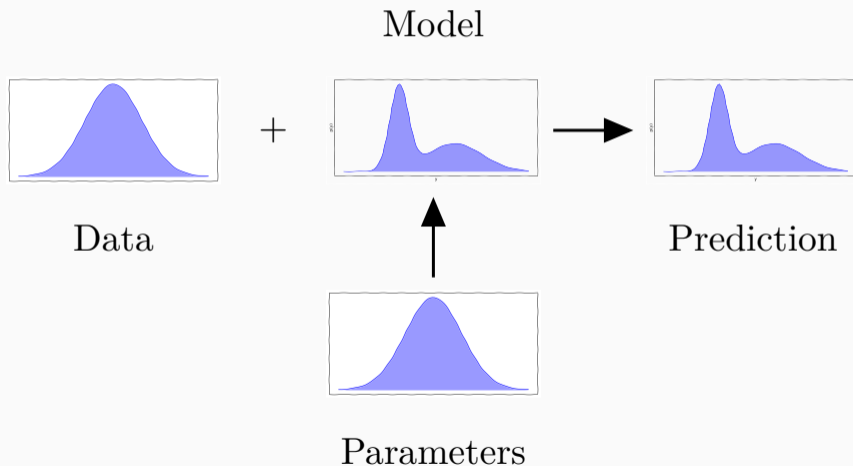


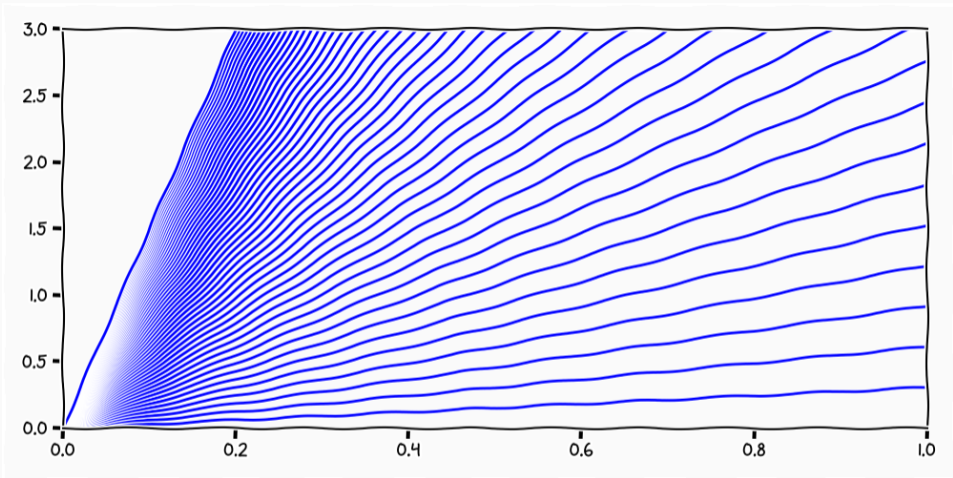




$$\mathbf{y} = \underbrace{f_N}_{w_N} \circ \underbrace{f_{N-1}}_{w_{N-1}} \circ \cdots \circ \underbrace{f_1}_{w_1} \circ \underbrace{f_0}_{w_0} (\mathbf{x})$$

Parametric Knowledge (Uncertainty) Propagation





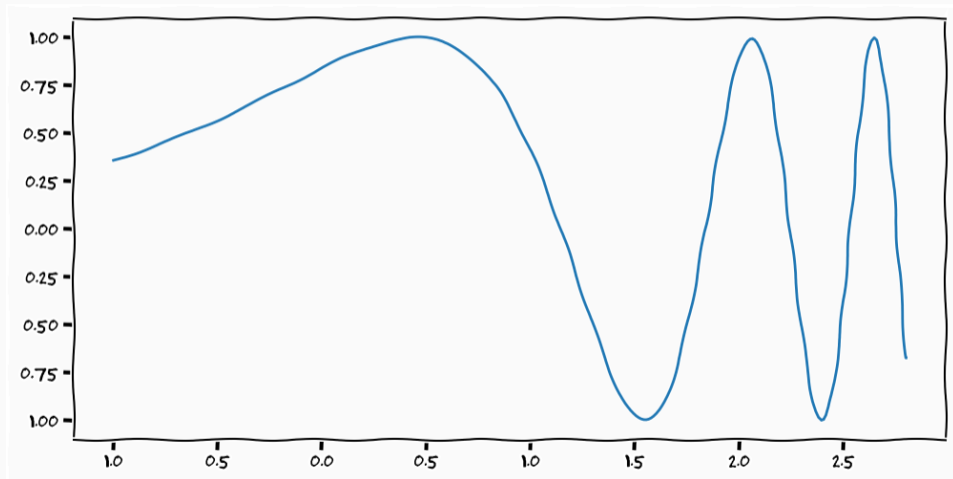
$$\mathbf{y} = \underbrace{f_N}_{w_N} \circ \underbrace{f_{N-1}}_{w_{N-1}} \circ \cdots \circ \underbrace{f_1}_{w_1} \circ \underbrace{f_0}_{w_0} (\mathbf{x})$$

When do I want Composite Functions

$$y = f_k \circ f_{k-1} \circ \cdots \circ f_1(x)$$

1. My generative process is composite
 - my prior knowledge is composite
2. I want to "re-parametrise" my kernel in a learning setting
 - i have knowledge of the re-parametrisation

Because we lack "models"?



Diff Levels of Abstraction

- Hierarchical Learning

- Natural progression from low level to high level structure as seen in natural complexity
- Easier to monitor what is being learnt and to guide the machine to better subspaces
- A good lower level representation can be used for many distinct tasks

Feature representation



3rd layer
"Objects"



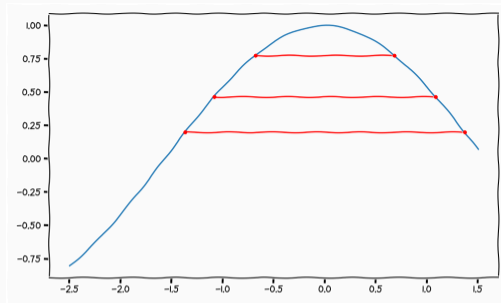
2nd layer
"Object parts"



1st layer
"Edges"



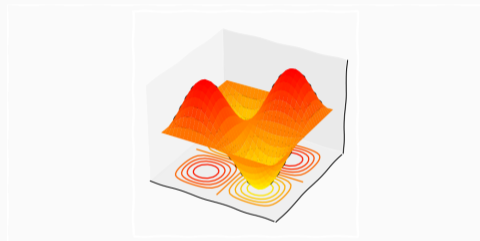
Pixels



$$y = f_k(f_{k-1}(\dots f_0(x))) = f_k \circ f_{k-1} \circ \dots \circ f_1(x)$$

$$\text{Kern}(f_1) \subseteq \text{Kern}(f_{k-1} \circ \dots \circ f_2 \circ f_1) \subseteq \text{Kern}(f_k \circ f_{k-1} \circ \dots \circ f_2 \circ f_1)$$

$$\text{Im}(f_k \circ f_{k-1} \circ \dots \circ f_2 \circ f_1) \subseteq \text{Im}(f_k \circ f_{k-1} \circ \dots \circ f_2) \subseteq \dots \subseteq \text{Im}(f_k)$$

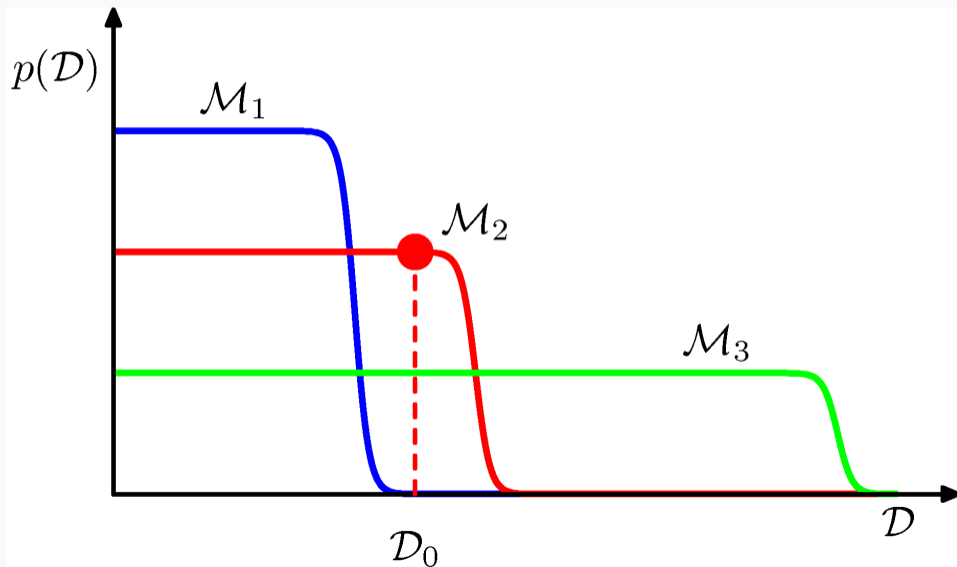




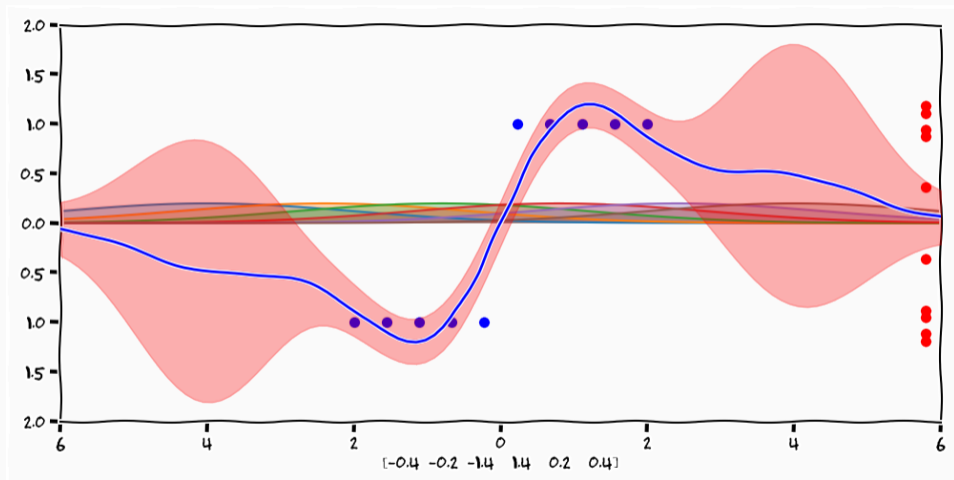


$$p(\mathbf{y}) = \int p(\mathbf{y} \mid w_N) p(w_N \mid w_{N-1}) \cdots p(w_1 \mid w_0) p(w_0) dw_N dw_{N-1} dw_1 dw_0$$

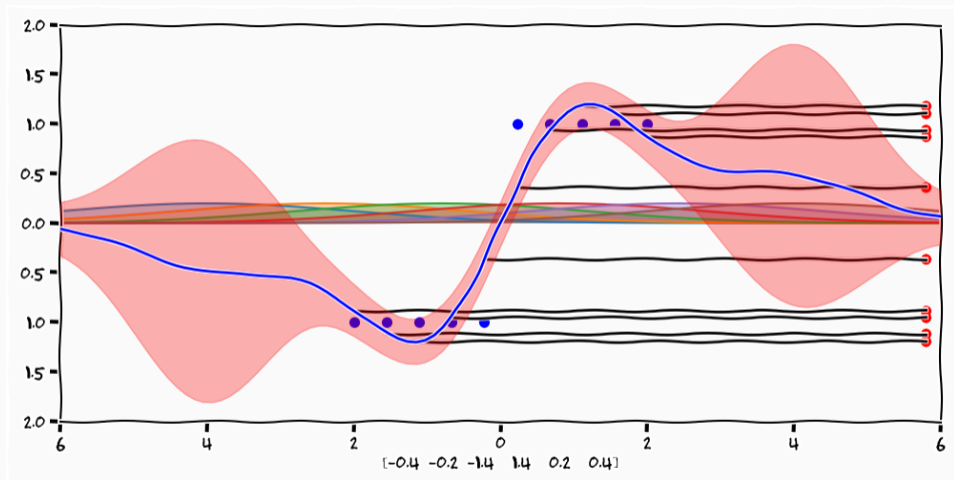
MacKay plot



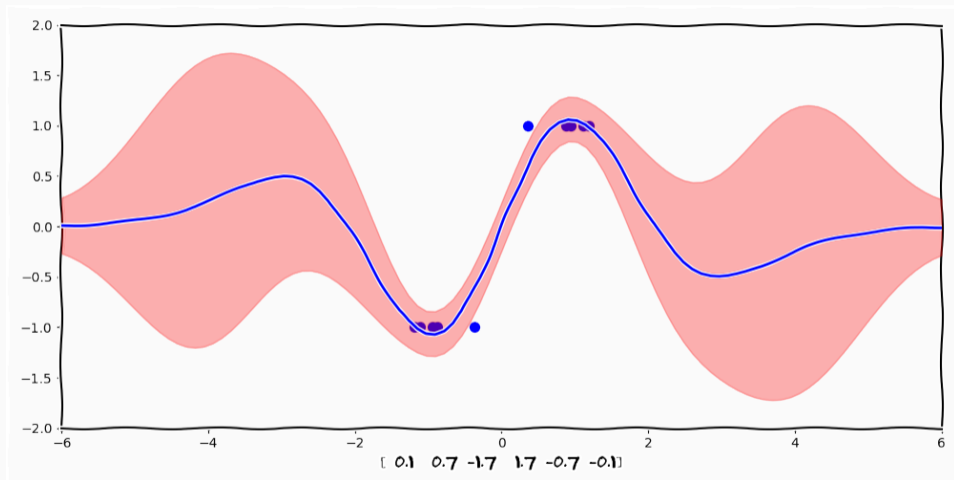
Composite Functions



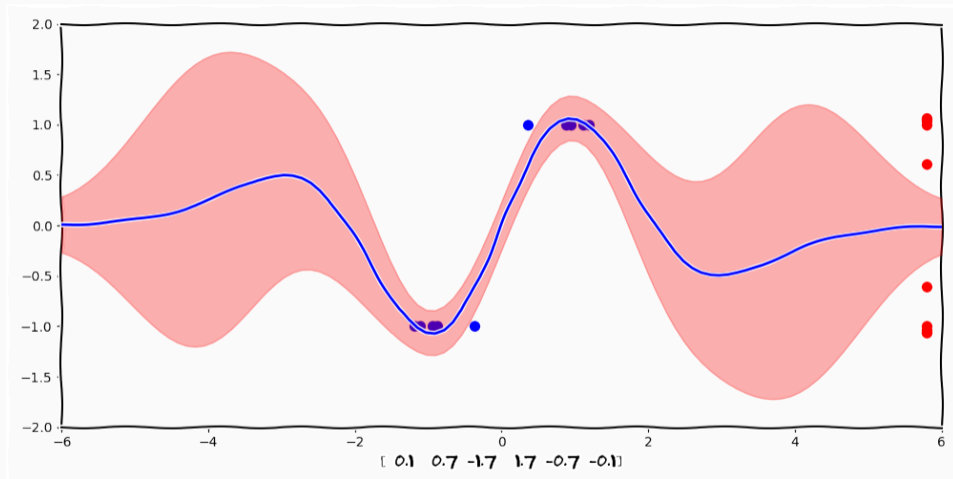
Composite Functions



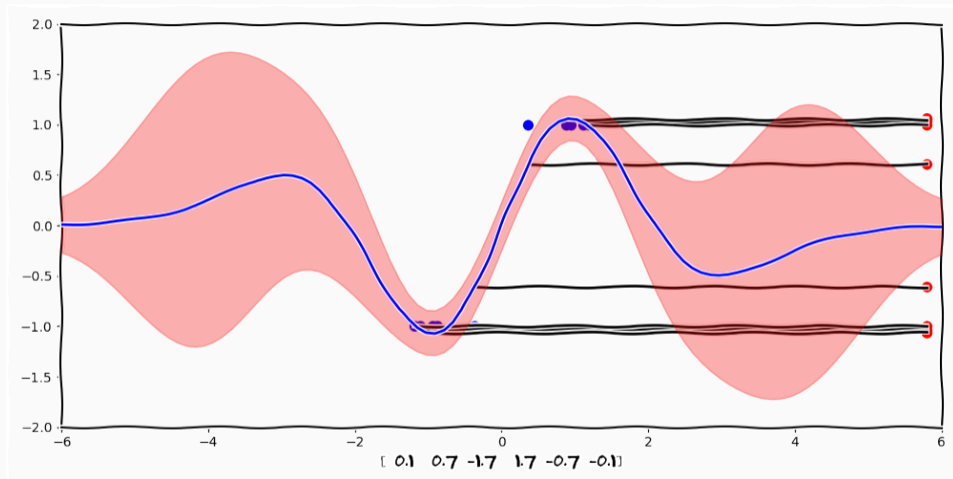
Composite Functions



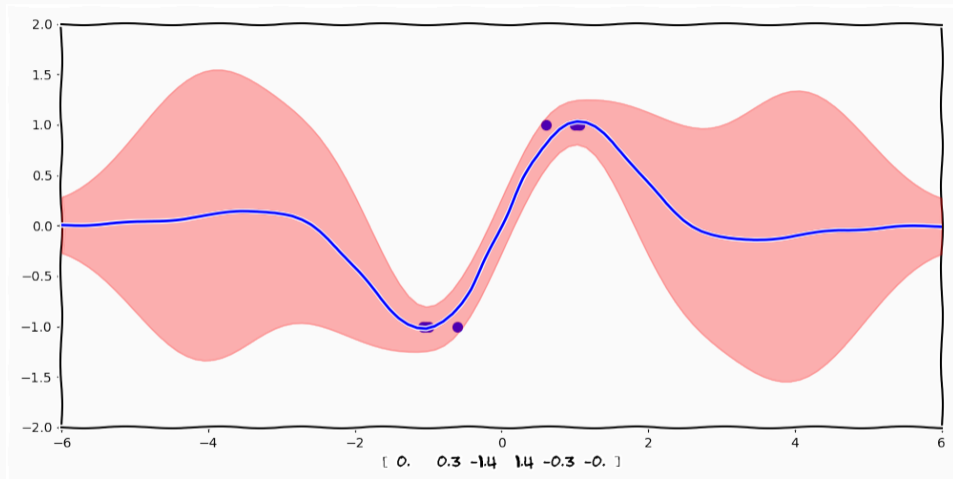
Composite Functions



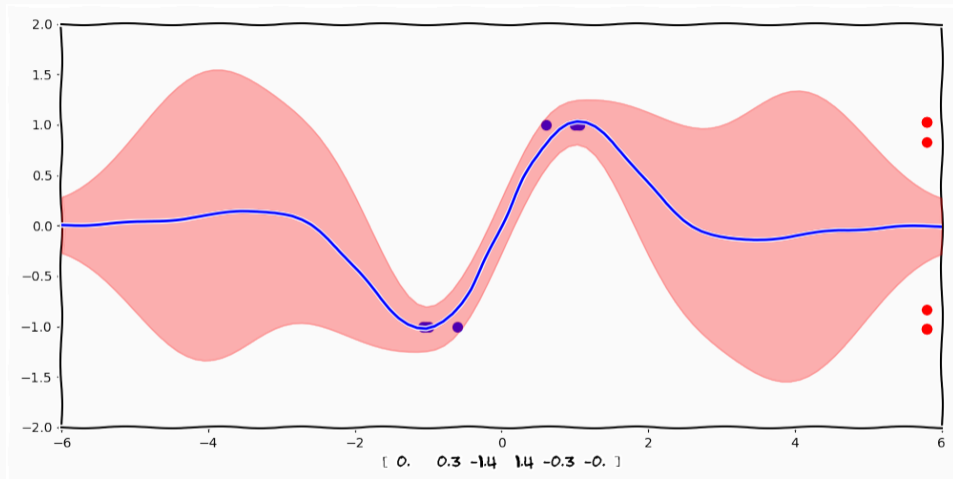
Composite Functions



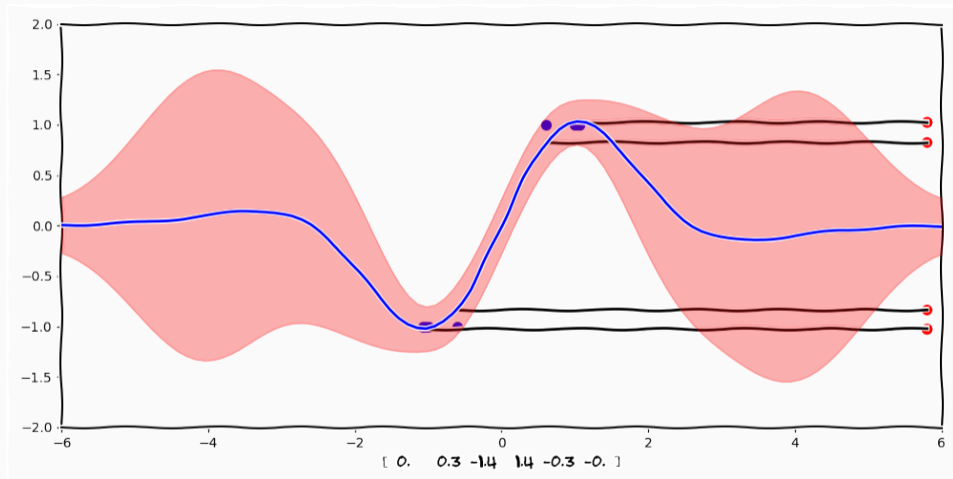
Composite Functions



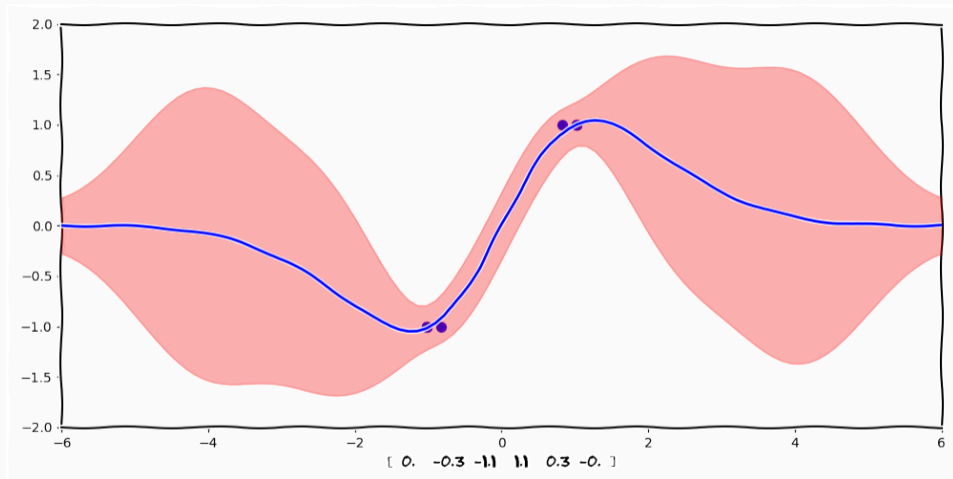
Composite Functions



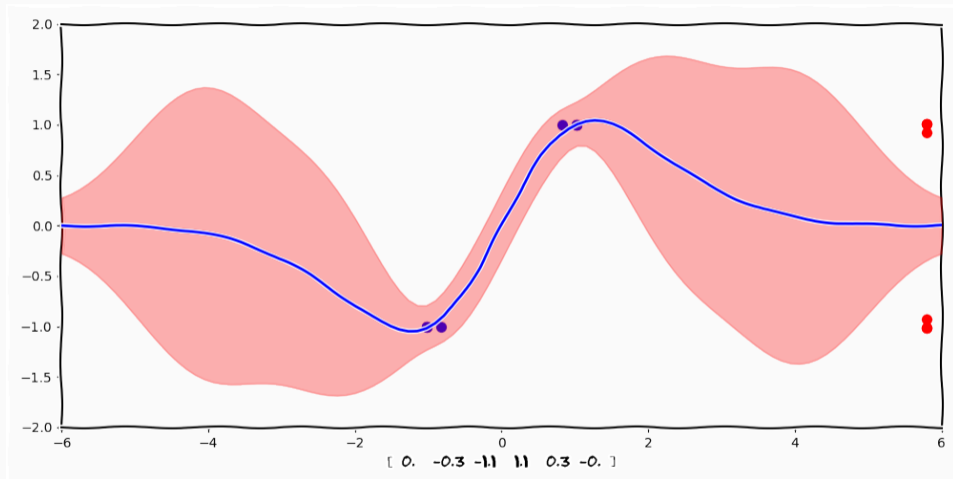
Composite Functions



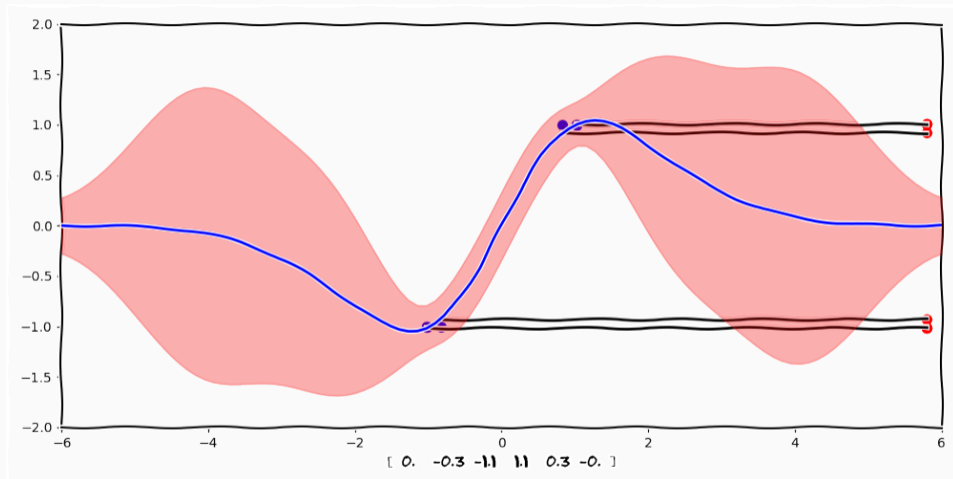
Composite Functions



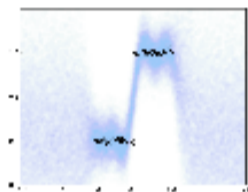
Composite Functions



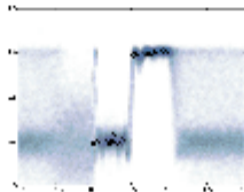
Composite Functions



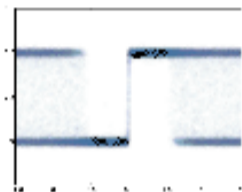
The Final Composition



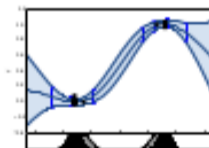
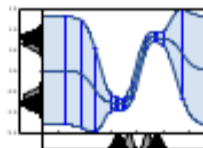
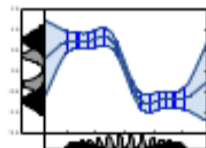
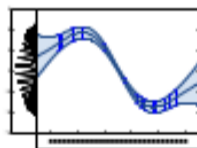
(a) GP



(b) 2 layers

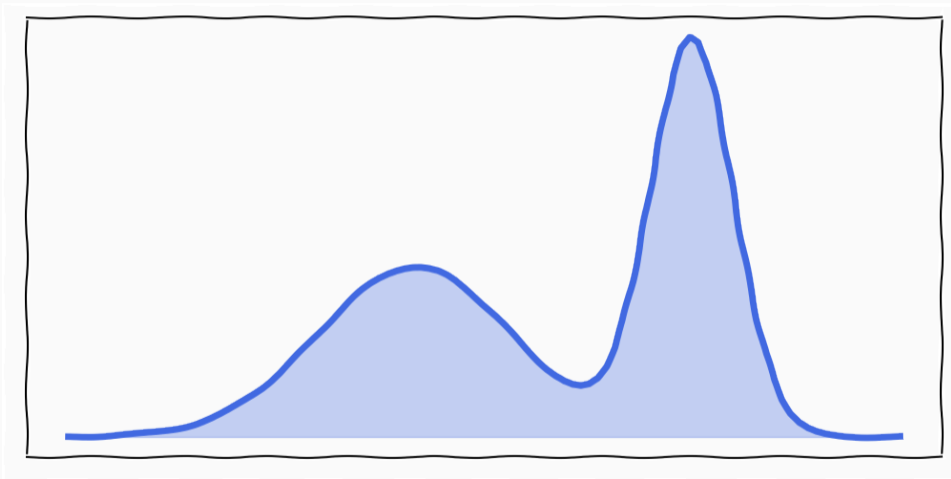


(c) 4 layers

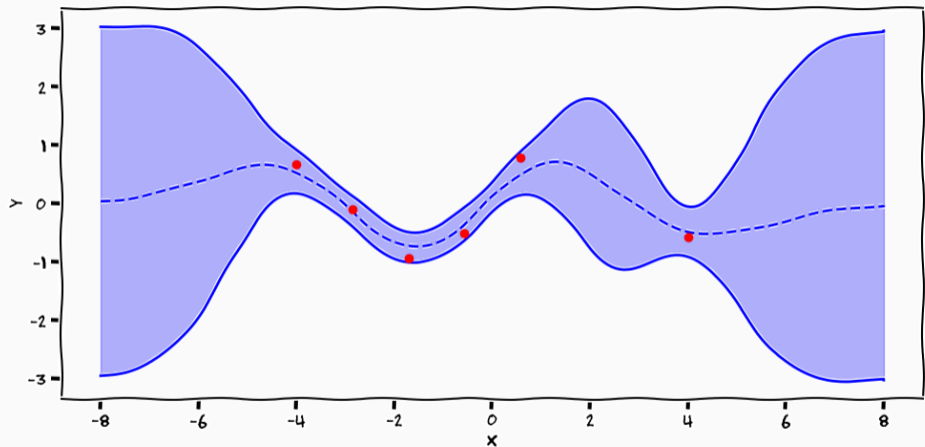


(d) Hidden spaces for 4 layer model

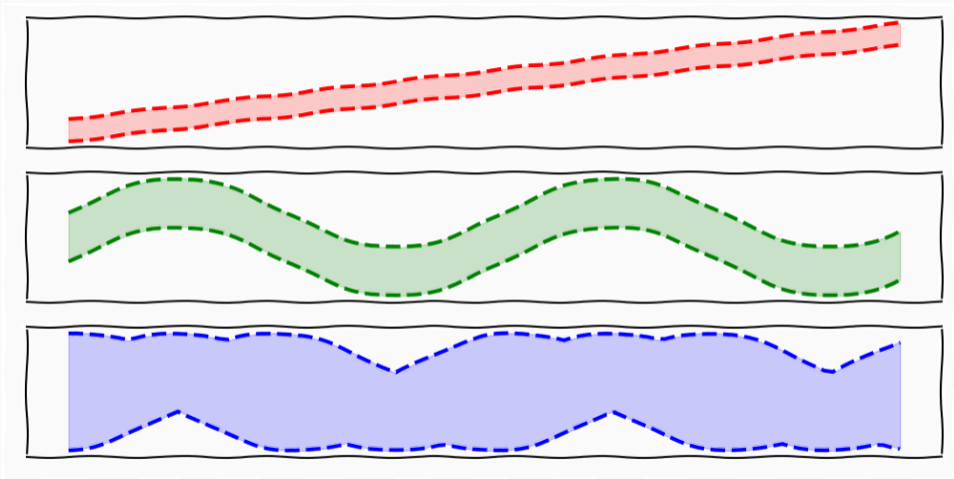
Remember why we did this in the first place



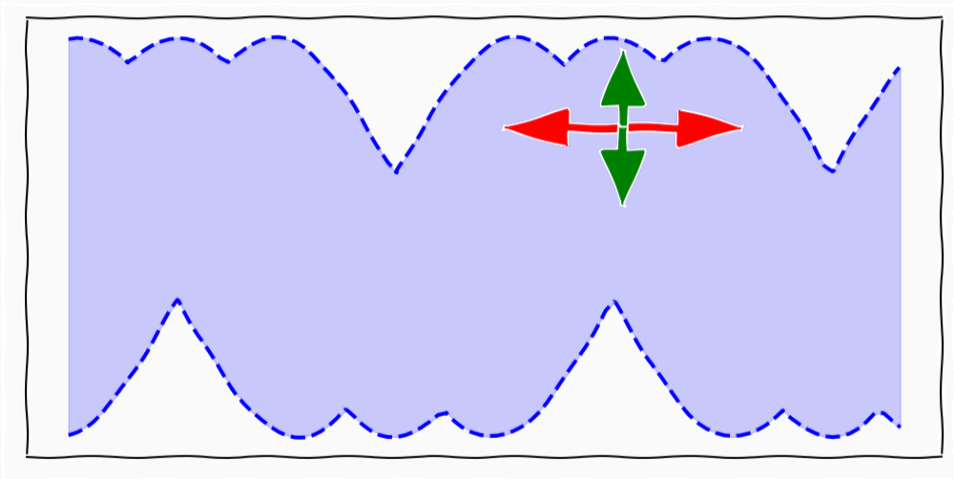
These damn plots



It gets worse



It gets even worse

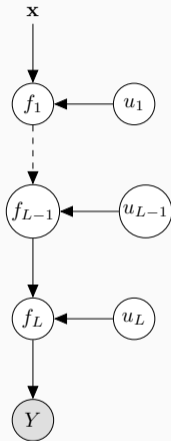


- Sufficient statistics

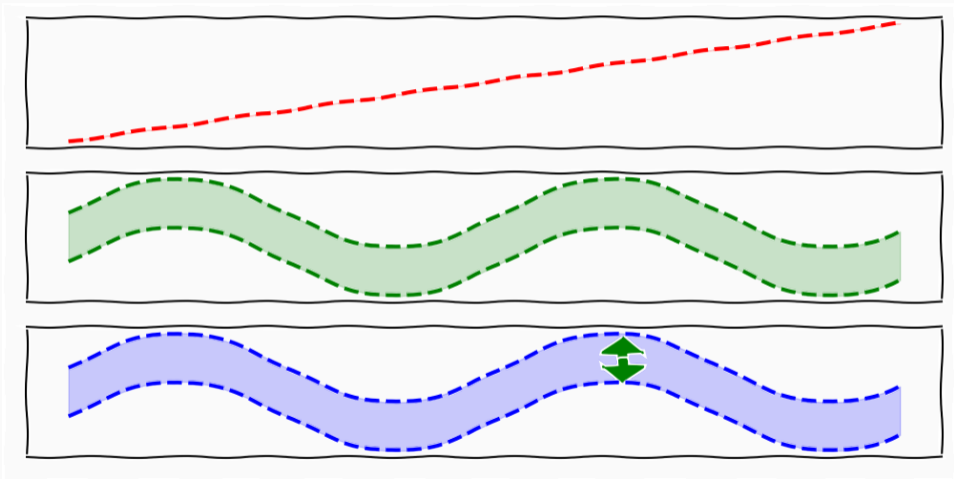
$$\begin{aligned} q(\mathbf{F})q(\mathbf{U})q(\mathbf{X}) &= p(\mathbf{F}|\mathbf{Y}, \mathbf{U}, \mathbf{X}, \mathbf{Z})q(\mathbf{U})q(\mathbf{X}) \\ &= p(\mathbf{F}|\mathbf{U}, \mathbf{X}, \mathbf{Z})q(\mathbf{U})q(\mathbf{X}) \end{aligned}$$

- Mean-Field

$$q(\mathbf{U}) = \prod_i^L q(\mathbf{U}_i)$$



The effect



What have we lost

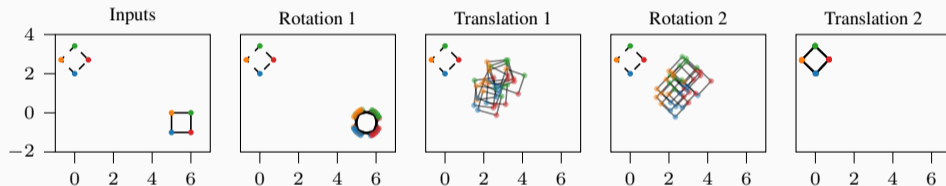
- Our priors are not reflected correctly
 - → we cannot interpret the results
- No intermediate uncertainties
 - → we cannot do sequential decision making

What have we lost

- Our priors are not reflected correctly
 - → we cannot interpret the results
- No intermediate uncertainties
 - → we cannot do sequential decision making
- We are performing a massive computational overhead for very little use

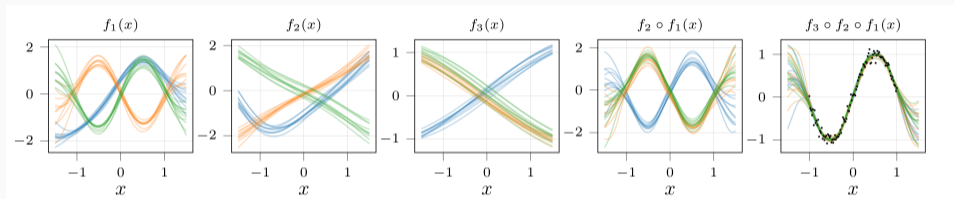
- Our priors are not reflected correctly
 - → we cannot interpret the results
- No intermediate uncertainties
 - → we cannot do sequential decision making
- We are performing a massive computational overhead for very little use
- *"...throwing out the baby with the bathwater..."*

What we really want¹



¹Ustyuzhaninov et al., 2020

What we really want²



★

²Ustyuzhaninov et al., 2020

- The community have tried Bayesian principles for composite functions for a long time MacKay, [1991](#)

- The community have tried Bayesian principles for composite functions for a long time MacKay, [1991](#)
- Empirically performance of Bayesian inference of composite functions is not impressive

- The community have tried Bayesian principles for composite functions for a long time MacKay, [1991](#)
- Empirically performance of Bayesian inference of composite functions is not impressive
- It is not just parametric models, non-parametric composite models doesn't work either

- The community have tried Bayesian principles for composite functions for a long time MacKay, 1991
- Empirically performance of Bayesian inference of composite functions is not impressive
- It is not just parametric models, non-parametric composite models doesn't work either
- *Why?*

- Roy, H., Miani, M., Ek, C. H., Hennig, P., Pförtner, M., Tatzel, L., & Hauberg, S., (2024). **Reparameterization invariance in approximate bayesian inference**. In Advances in Neural Information Processing Systems (NeurIPS)
- Fadel, S., Roy, H., Krämer, N., Zainchkovskyy, Y., Syrota, S., Mahou, A., Ek, C. H., Hauberg, S. (2025). **Deep variational inference with stochastic projections**. In Submission

Statistical Models of Composite Functions

$$\mathbf{y} = f(\mathbf{x}) = \underbrace{f_N}_{w_N} \circ \underbrace{f_{N-1}}_{w_{N-1}} \circ \cdots \circ \underbrace{f_1}_{w_1} \circ \underbrace{f_0}_{w_0}(\mathbf{x})$$

$$p(\mathbf{w} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \alpha \mathbf{I})$$

$$p(\mathbf{w} \mid \mathcal{D}) = \frac{1}{Z} p(\mathcal{D} \mid \mathbf{w}) p(\mathbf{w}) = \frac{1}{Z} \exp(-\mathcal{L}(\mathcal{D}; \mathbf{w}))$$
$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \mathcal{L}(\mathcal{D}; \mathbf{w})$$

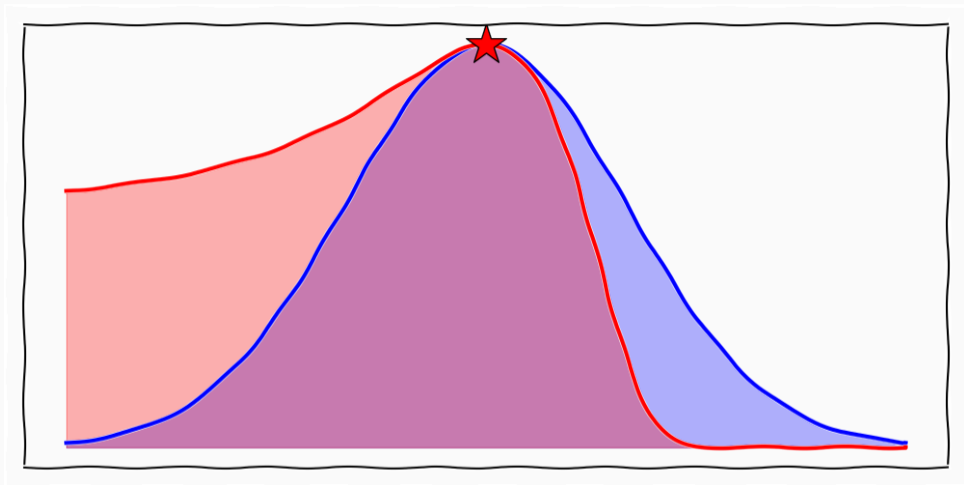
$$\begin{aligned}\mathcal{L}(\mathcal{D}; \mathbf{w}) &\approx \mathcal{L}(\mathcal{D}; \hat{\mathbf{w}}) + \nabla \mathcal{L}(\mathcal{D}; \mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}} (\mathbf{w} - \hat{\mathbf{w}}) \\ &\quad + \frac{1}{2}(\mathbf{w} - \hat{\mathbf{w}})^T \nabla^2 \mathcal{L}(\mathcal{D}; \mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}} (\mathbf{w} - \hat{\mathbf{w}})\end{aligned}$$

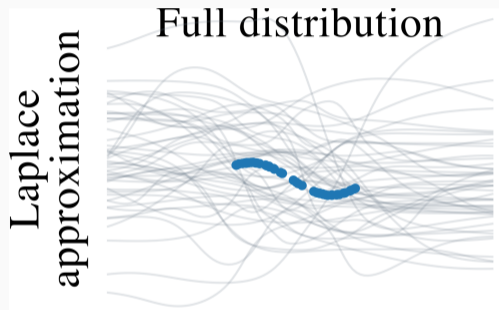
$$\begin{aligned}\mathcal{L}(\mathcal{D}; \mathbf{w}) &\approx \mathcal{L}(\mathcal{D}; \hat{\mathbf{w}}) + \nabla \mathcal{L}(\mathcal{D}; \mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}} (\mathbf{w} - \hat{\mathbf{w}}) \\ &\quad + \frac{1}{2}(\mathbf{w} - \hat{\mathbf{w}})^{\mathrm{T}} \nabla^2 \mathcal{L}(\mathcal{D}; \mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}} (\mathbf{w} - \hat{\mathbf{w}}) \\ &= \mathcal{L}(\mathcal{D}; \hat{\mathbf{w}}) + \frac{1}{2}(\mathbf{w} - \hat{\mathbf{w}})^{\mathrm{T}} \nabla^2 \mathcal{L}(\mathcal{D}; \mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}} (\mathbf{w} - \hat{\mathbf{w}})\end{aligned}$$

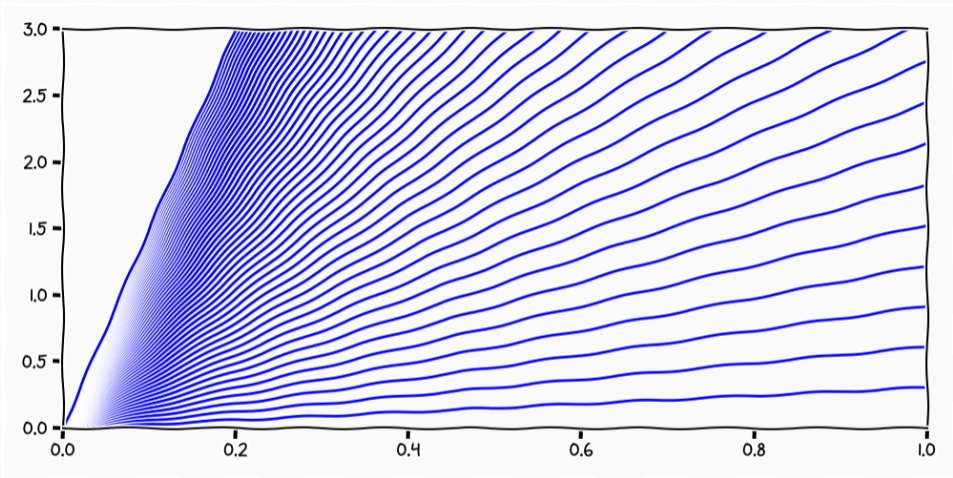
$$p(\mathbf{w} \mid \mathcal{D}) = \frac{1}{Z} p(\mathcal{D} \mid \mathbf{w}) p(\mathbf{w}) = \frac{1}{Z} \exp(-\mathcal{L}(\mathcal{D}; \mathbf{w}))$$

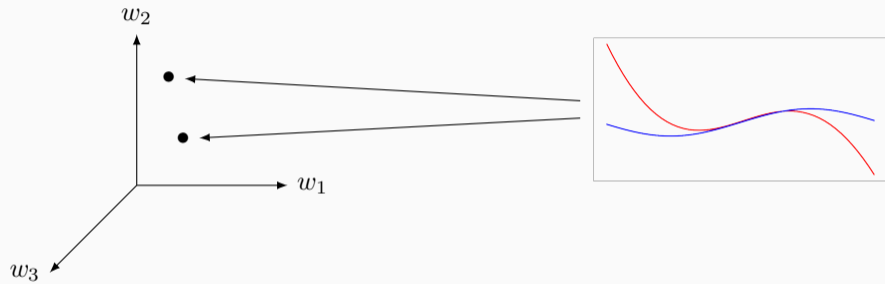
$$\begin{aligned} p(\mathbf{w} \mid \mathcal{D}) &= \frac{1}{Z} p(\mathcal{D} \mid \mathbf{w}) p(\mathbf{w}) = \frac{1}{Z} \exp(-\mathcal{L}(\mathcal{D}; \mathbf{w})) \\ &= \frac{1}{Z} \exp(-\mathcal{L}(\mathcal{D}; \hat{\mathbf{w}})) \exp\left(\frac{1}{2}(\mathbf{w} - \hat{\mathbf{w}})^T \nabla^2 \mathcal{L}(\mathcal{D}; \mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}} (\mathbf{w} - \hat{\mathbf{w}})\right) \end{aligned}$$

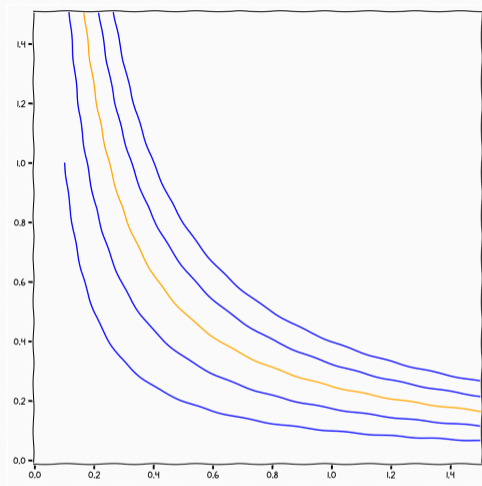
$$\begin{aligned} p(\mathbf{w} \mid \mathcal{D}) &= \frac{1}{Z} p(\mathcal{D} \mid \mathbf{w}) p(\mathbf{w}) = \frac{1}{Z} \exp(-\mathcal{L}(\mathcal{D}; \mathbf{w})) \\ &= \frac{1}{Z} \exp(-\mathcal{L}(\mathcal{D}; \hat{\mathbf{w}})) \exp\left(\frac{1}{2}(\mathbf{w} - \hat{\mathbf{w}})^T \nabla^2 \mathcal{L}(\mathcal{D}; \mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}} (\mathbf{w} - \hat{\mathbf{w}})\right) \\ &= \mathcal{N}\left(\mathbf{w} \mid \hat{\mathbf{w}}, (\nabla^2 \mathcal{L}(\mathcal{D}; \mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}})^{-1}\right) \end{aligned}$$



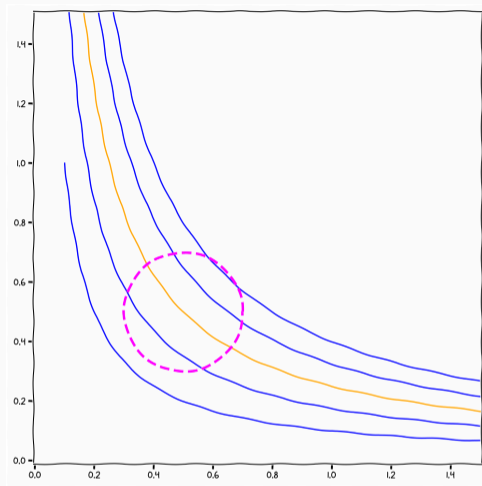




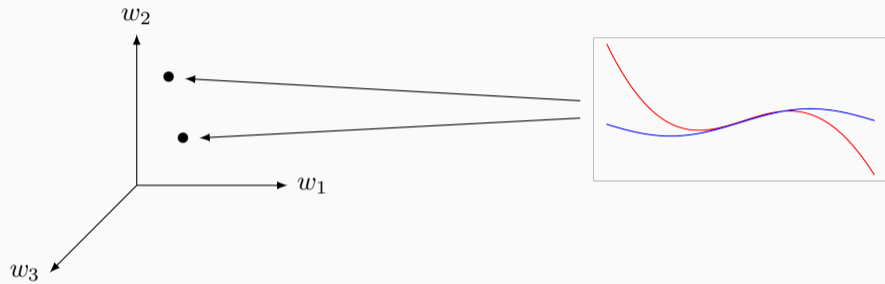


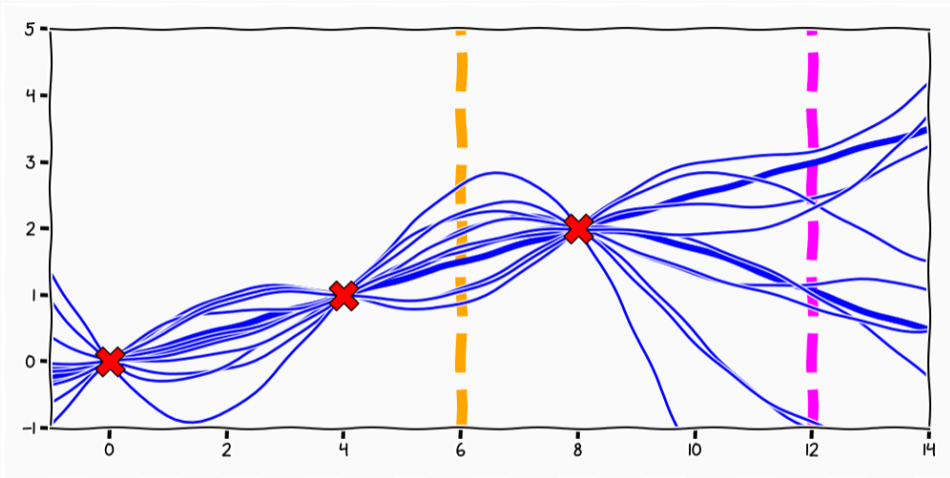


$$f(x) = w_1 \cdot w_2 \cdot x = (\alpha \cdot w_1) \cdot \left(\frac{1}{\alpha} \cdot w_2 \right) \cdot x$$



$$f(x) = w_1 \cdot w_2 \cdot x = (\alpha \cdot w_1) \cdot \left(\frac{1}{\alpha} \cdot w_2 \right) \cdot x$$





Definition (x-reparametrisations)

Given a datapoint $\mathbf{x} \in \mathbb{R}^I$, for any $\mathbf{w} \in \mathbb{R}^D$ we define the \mathbf{x} -reparameterizations as the set

$\mathcal{R}_{\mathbf{x}}^f(\mathbf{w}) = \{\mathbf{w}' \text{ such that } f(\mathbf{w}', \mathbf{x}) = f(\mathbf{w}, \mathbf{x})\}$. Consistently, given a collection of points $\mathcal{X} \subseteq \mathbb{R}^I$, we call the intersection

$\mathcal{R}_{\mathcal{X}}^f(\mathbf{w}) = \bigcap_{\mathbf{x} \in \mathcal{X}} \mathcal{R}_{\mathbf{x}}^f(\mathbf{w})$ \mathcal{X} -reparameterizations.

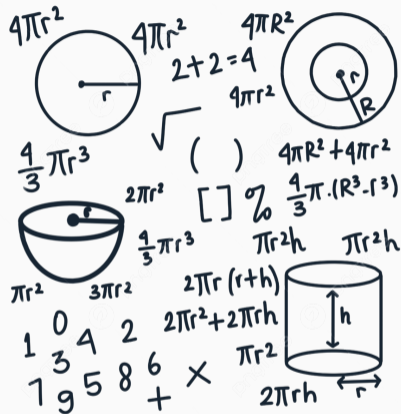
- We define the relation \sim over \mathbb{R}^D as $\mathbf{w} \sim \mathbf{w}'$ if $\mathbf{w}' \in \bar{\mathcal{R}}_{\mathcal{X}}^f(\mathbf{w})$.

The effective parameter quotient group

- We define the relation \sim over \mathbb{R}^D as $\mathbf{w} \sim \mathbf{w}'$ if $\mathbf{w}' \in \bar{\mathcal{R}}_{\mathcal{X}}^f(\mathbf{w})$.
- Quotient space of effective parameters $\mathcal{P} = \mathbb{R}^D / \sim$

The effective parameter quotient group

- We define the relation \sim over \mathbb{R}^D as $\mathbf{w} \sim \mathbf{w}'$ if $\mathbf{w}' \in \bar{\mathcal{R}}_{\mathcal{X}}^f(\mathbf{w})$.
- Quotient space of effective parameters $\mathcal{P} = \mathbb{R}^D / \sim$
- $[\mathbf{w}_1], [\mathbf{w}_2] \in \mathcal{P}$ are the same point if and only if $\mathbf{w}_1 \sim \mathbf{w}_2$.



$$\text{dist}(\mathbf{w}_1, \mathbf{w}_2) = 0 \leftrightarrow \mathbf{w}_1 \sim \mathbf{w}_2$$

$$\text{dist}(\mathbf{w}_1, \mathbf{w}_2) = 0 \Leftrightarrow \mathbf{w}_1 \sim \mathbf{w}_2$$

$$\text{dist}^2(\mathbf{w}, \mathbf{w} + \boldsymbol{\epsilon}) = \sum_{n=1}^N ||f(\mathbf{w}, \mathbf{x}_n) - f(\mathbf{w} + \boldsymbol{\epsilon}, \mathbf{x}_n)||^2$$

$$\text{dist}(\mathbf{w}_1, \mathbf{w}_2) = 0 \leftrightarrow \mathbf{w}_1 \sim \mathbf{w}_2$$

$$\begin{aligned} \text{dist}^2(\mathbf{w}, \mathbf{w} + \boldsymbol{\epsilon}) &= \sum_{n=1}^N ||f(\mathbf{w}, \mathbf{x}_n) - f(\mathbf{w} + \boldsymbol{\epsilon}, \mathbf{x}_n)||^2 \\ &= \sum_{n=1}^N ||f(\mathbf{w}, \mathbf{x}_n) - f(\mathbf{w}, \mathbf{x}_n) - \nabla f(\mathbf{w}, \mathbf{x}_n)\boldsymbol{\epsilon}||^2 \end{aligned}$$

$$\text{dist}(\mathbf{w}_1, \mathbf{w}_2) = 0 \leftrightarrow \mathbf{w}_1 \sim \mathbf{w}_2$$

$$\begin{aligned} \text{dist}^2(\mathbf{w}, \mathbf{w} + \boldsymbol{\epsilon}) &= \sum_{n=1}^N \|f(\mathbf{w}, \mathbf{x}_n) - f(\mathbf{w} + \boldsymbol{\epsilon}, \mathbf{x}_n)\|^2 \\ &= \sum_{n=1}^N \|f(\mathbf{w}, \mathbf{x}_n) - f(\mathbf{w}, \mathbf{x}_n) - \nabla f(\mathbf{w}, \mathbf{x}_n)\boldsymbol{\epsilon}\|^2 \\ &= \boldsymbol{\epsilon}^T \mathbf{J}^T \mathbf{J} \boldsymbol{\epsilon} \end{aligned}$$

$$q(\mathbf{w}) = \mathcal{N} \left(\mathbf{w} \mid \hat{\mathbf{w}}, \left(\underbrace{\mathbf{J}^T \nabla^2 \mathcal{L}(\mathcal{D}; \mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}} \mathbf{J}}_{\mathbf{GGN}_{\mathbf{w}}} \right)^{-1} \right)$$

- This is also known as the Generalised Gauss Newton Approximation

$$q(\mathbf{w}) = \mathcal{N} \left(\mathbf{w} \mid \hat{\mathbf{w}}, \left(\underbrace{\mathbf{J}^T \nabla^2 \mathcal{L}(\mathcal{D}; \mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}} \mathbf{J}}_{\mathbf{GGN}_{\mathbf{w}}} \right)^{-1} \right)$$

- This is also known as the Generalised Gauss Newton Approximation
- Linearised Laplace Approximation Immer et al., [2021](#)

$$q(\mathbf{w}) = \mathcal{N} \left(\mathbf{w} \mid \hat{\mathbf{w}}, \left(\underbrace{\mathbf{J}^T \nabla^2 \mathcal{L}(\mathcal{D}; \mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}} \mathbf{J}}_{\mathbf{GGN}_{\mathbf{w}}} \right)^{-1} \right)$$

- This is also known as the Generalised Gauss Newton Approximation
- Linearised Laplace Approximation Immer et al., [2021](#)
- Interpreted as a Riemannian metric it is called the Fisher-Rao metric

$$q(\mathbf{w}) = \mathcal{N} \left(\mathbf{w} \mid \hat{\mathbf{w}}, \left(\underbrace{\mathbf{J}^T \nabla^2 \mathcal{L}(\mathcal{D}; \mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}} \mathbf{J}}_{\mathbf{GGN}_{\mathbf{w}}} \right)^{-1} \right)$$

- This is also known as the Generalised Gauss Newton Approximation
- Linearised Laplace Approximation Immer et al., [2021](#)
- Interpreted as a Riemannian metric it is called the Fisher-Rao metric
- It is a *pseudo-metric*

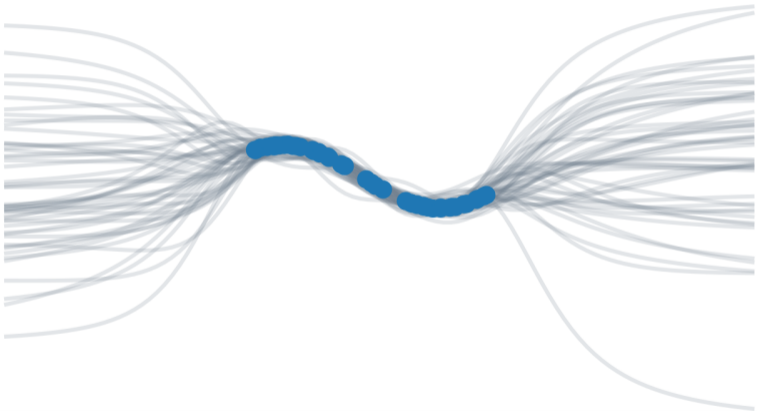
Theorem (Topological Equivalence)

The pseudo-Riemannian manifold obtained with the pullback pseudo-metric $(\mathbb{R}^D, GGN_{\mathbf{w}})$ is homeomorphic to the quotient group $(\mathcal{P}, d_{\mathcal{P}})$

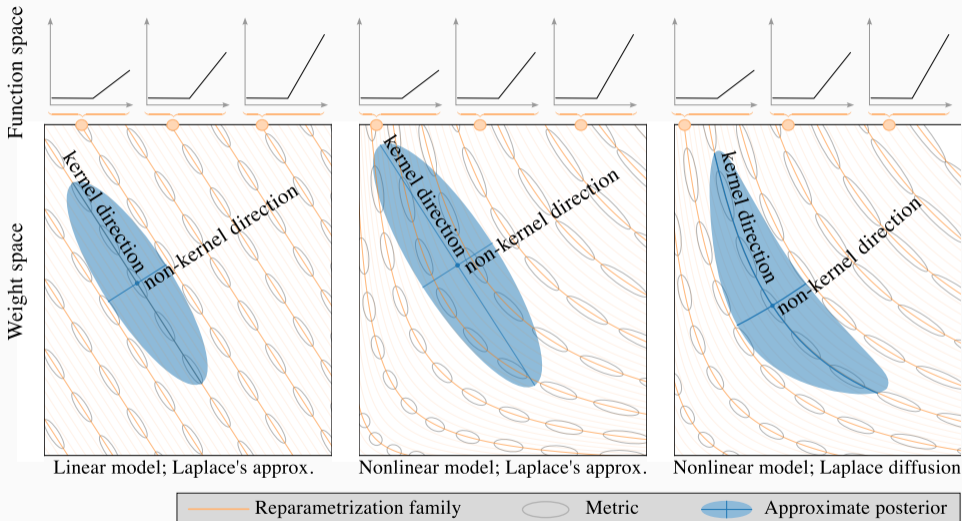
For any $\mathbf{w}_0, \mathbf{w}_1 \in \mathbf{R}^D$ it holds

$$d_{f^*H}(\mathbf{w}_0, \mathbf{w}_1) = 0 \quad \Longleftrightarrow \quad [\mathbf{w}_0] = [\mathbf{w}_1] \in \mathcal{P}$$

Linearized Laplace



Intuition



- Linear Function

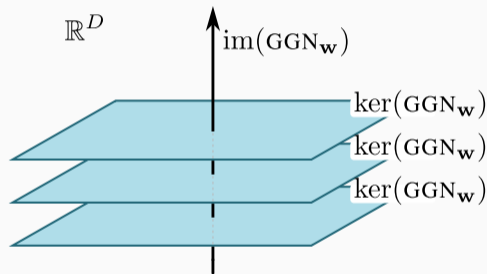
$$f(\mathbf{w}) = \mathbf{A}\mathbf{w} + b$$

$$g : \mathbb{R}^D \rightarrow \mathbb{R}^D$$

$$\text{s.t. } \mathbf{A}(g(\mathbf{w}) - \mathbf{w}) = \mathbf{0}$$

- Nullspace of \mathbf{A}

$$(g(\mathbf{w}) - \mathbf{w}) \in \ker(\mathbf{A})$$



$$f_w(x) \approx f_{\hat{w}} + \mathbf{J}_w(x)(w - \hat{w})$$

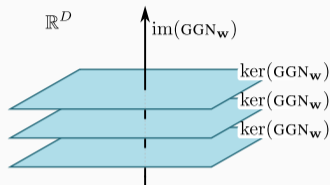
- re-parametrisations are characterised by the kern(\mathbf{J}_w)
- By construction

$$\text{kern}(\mathbf{J}_w) = \text{kern}(\mathbf{J}_w^T \mathbf{J}_w)$$

- Neural Tangent Kernel Jacot et al., [2018](#)

$$\text{NTK} = \mathbf{J}_w \mathbf{J}_w^T$$

Orthogonal Subspaces



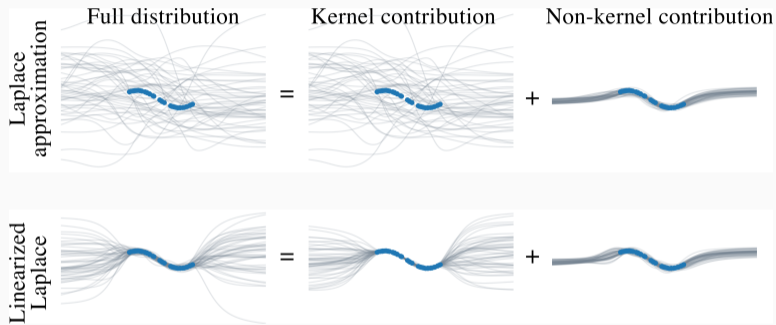
$$\text{im}(\text{GGN}_{\mathbf{w}}) \oplus \text{kern}(\text{GGN}_{\mathbf{w}}) = \mathbb{R}^D$$

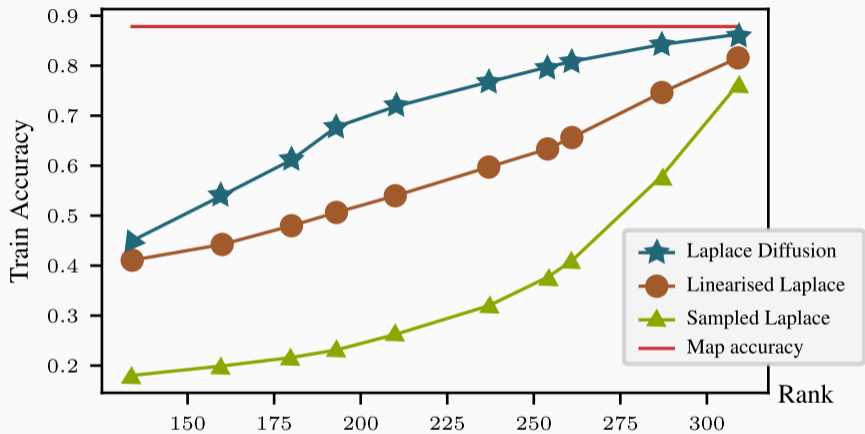
- $\text{im}(\text{GGN}_{\mathbf{w}})$ - spans the *effective* parameters of the model
- $\text{kern}(\text{GGN}_{\mathbf{w}})$ - parameters leading to the same function

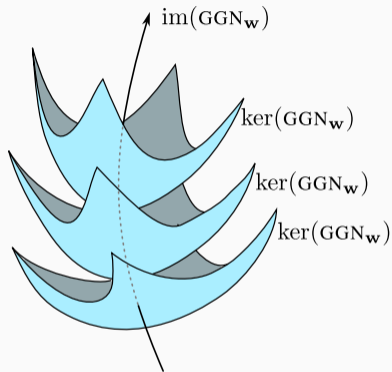
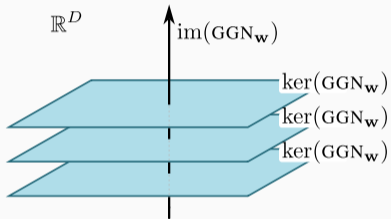
$$\Sigma = \left(\begin{bmatrix} U_1 \\ U_2 \end{bmatrix}^T \begin{bmatrix} \tilde{\Lambda} & | & 0 \\ \hline 0 & | & 0 \end{bmatrix} \begin{bmatrix} U_1 \\ U_2 \end{bmatrix} + \alpha I \right)^{-1} = U_1^T (\tilde{\Lambda} + \alpha I_k)^{-1} U_1 + \alpha^{-1} U_2^T U_2.$$

- Decomposition of parameter space

$$\mathbf{w} = \hat{\mathbf{w}} + \mathbf{w}_{\text{ker}} + \mathbf{w}_{\text{im}}$$







- Riemannian Diffusion

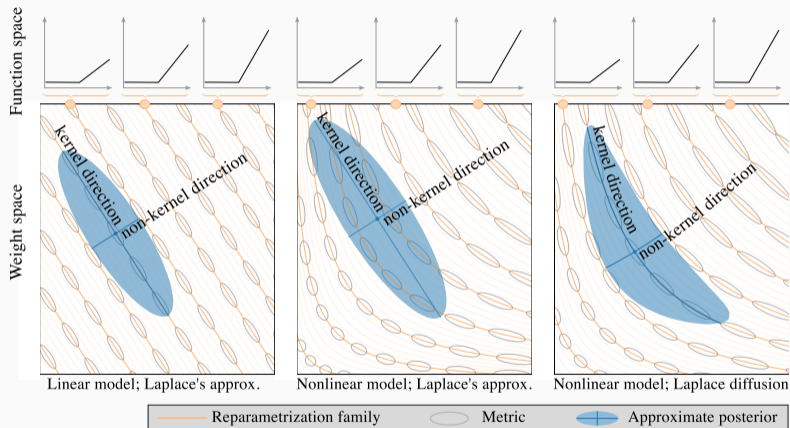
$$\mathbf{w} = \sqrt{2\tau} G(\mathbf{w})^{-\frac{1}{2}} W + \tau \Gamma t$$

$$\text{where } \Gamma_i(\mathbf{w}) = \sum_{j=1}^D \frac{\partial}{\partial \mathbf{w}_j} (G(\mathbf{w})^{-1})_{ij}.$$

- Update rule

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \sqrt{2h_t} G(\mathbf{w}_t)^{-\frac{1}{2}} \epsilon$$

Intuition



- Diffusion in Linearised Laplace

$$(\mathbb{R}^D, \mathbf{G}\mathbf{G}\mathbf{N}_{\hat{\mathbf{w}}} + \alpha\mathbf{I})$$

- Diffusion in Linearised Laplace

$$(\mathbb{R}^D, \mathbf{G}\mathbf{G}\mathbf{N}_{\hat{\mathbf{w}}} + \alpha\mathbf{I})$$

- Kernel Manifold

$$(\mathcal{P}_{\mathbf{w}}^{\perp}, \alpha\mathbf{I})$$

- Diffusion in Linearised Laplace

$$(\mathbb{R}^D, \mathbf{G}\mathbf{G}\mathbf{N}_{\hat{\mathbf{w}}} + \alpha\mathbf{I})$$

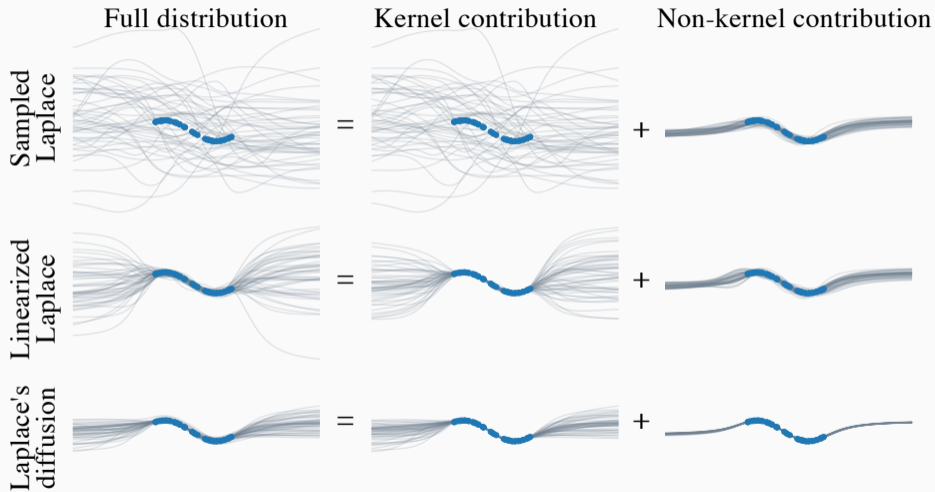
- Kernel Manifold

$$(\mathcal{P}_{\mathbf{w}}^{\perp}, \alpha\mathbf{I})$$

- Non-kernel Manifold

$$(\mathcal{P}_{\mathbf{w}}, \mathbf{G}\mathbf{G}\mathbf{N}_{\mathbf{w}}^{\perp})$$

Results



- we have characterised the geometry of reparametrisations

- we have characterised the geometry of reparametrisations
- the geometry explains why linearised laplace approximation works

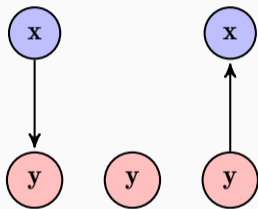
- we have characterised the geometry of reparametrisations
- the geometry explains why linearised laplace approximation works
- we have derived a simple random walk to sample from manifold

$$p(y) = \int p(y \mid \theta)p(\theta)\mathrm{d}\theta$$

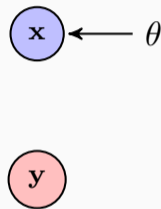
$$\begin{aligned} p(y) &= \int p(y \mid \theta) p(\theta) d\theta \\ &= \text{“}\tilde{p}(y_{\perp}) + \tilde{p}(y_{\neq})\text{”} \end{aligned}$$

Variational Inference

$$p(y) = \int p(y \mid x)p(x)dx$$



$$p(y) = \int_x p(y|x)p(x) = \frac{p(y|x)p(x)}{p(x|y)}$$



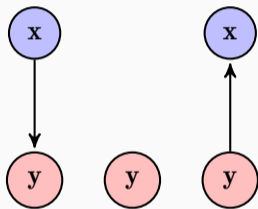
$$q_{\theta}(x) \approx p(x|y)$$

$$\log p(y) = \int q(x) \log \frac{p(x, y)}{p(x | y)} dx$$
$$=$$

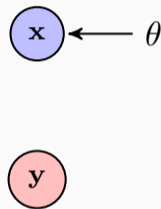
$$\begin{aligned}\log p(y) &= \int q(x) \log \frac{p(x, y)}{p(x | y)} dx \\ &= \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx + \int q(x) \log \frac{q(x)}{p(x|y)} dx\end{aligned}$$

$$\begin{aligned}\log p(y) &= \int q(x) \log \frac{p(x, y)}{p(x|y)} dx \\ &= \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx + \int q(x) \log \frac{q(x)}{p(x|y)} dx \\ &\geq - \int q(x) \log q(x) dx + \int q(x) \log p(x, y) dx\end{aligned}$$

- The Evidence Lower BOnd
- Tight if $q(x) = p(x|y)$



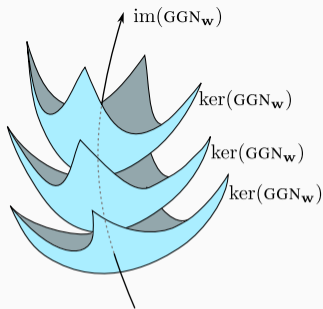
$$p(y) = \int_x p(y|x)p(x) = \frac{p(y|x)p(x)}{p(x|y)}$$



$$q_{\theta}(x) \approx p(x|y)$$

$$\mathcal{L}(q(x)) = \mathbb{E}_{q(x)} [\log p(x, y)] - H(q(x))$$

- We have to be able to compute an expectation over the joint distribution
- The second term should be trivial



$$\mathbf{J}^T \nabla^2 \mathcal{L}(\mathcal{D}; \mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}} \mathbf{J}$$

- Can we propose an approximate distribution that reflects the geometry of the parametrisation?

$$q(\theta) = \mathcal{N}(\theta | \hat{\theta}, \Sigma)$$
$$\Sigma = \sigma_{\text{ker}}^2 \mathbf{U} \mathbf{U}^T + \sigma^2 (\mathbb{I} - \mathbf{U} \mathbf{U}^T)$$

$$\mathcal{L} = \mathbb{E}_{\theta \sim q} [\log p(\mathbf{y}|\theta, \mathbf{x})] - \text{KL}(q(\theta) \| p(\theta))$$

Computing the Lower Bound using Stochastic Projections

$$\epsilon^{(s)} \sim \mathcal{N}(0, \mathbb{I}) \in \mathbb{R}^D,$$

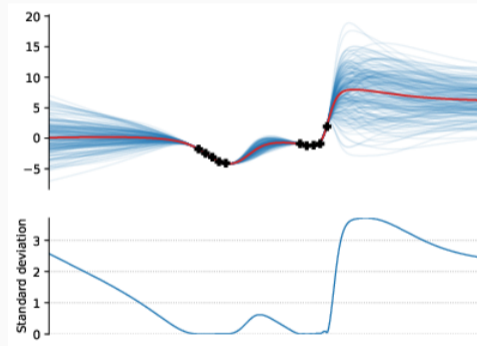
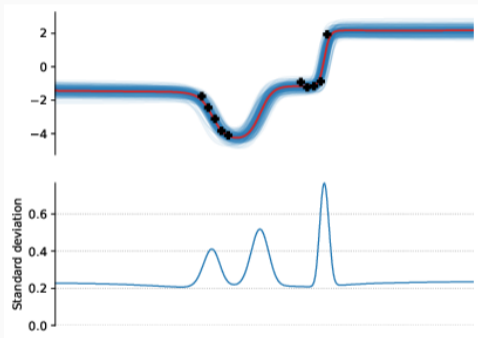
(projection onto kernel space) $\epsilon_{\text{ker}}^{(s)} = \mathbf{U}\mathbf{U}^T \epsilon^{(s)} \in \mathbb{R}^D,$

(image space is orthogonal) $\epsilon_{\text{im}}^{(s)} = (\mathbb{I} - \mathbf{U}\mathbf{U}^T) \epsilon^{(s)} = \epsilon^{(s)} - \epsilon_{\text{ker}}^{(s)} \in \mathbb{R}^D,$

$$\theta^{(s)} = \hat{\theta} + \sigma_{\text{ker}} \epsilon_{\text{ker}}^{(s)} + \sigma_{\text{im}} \epsilon_{\text{im}}^{(s)} \in \mathbb{R}^D,$$

$$\mathbb{E}_{\theta \sim q} [\log p(\theta, \mathbf{x})] \approx \frac{1}{S} \sum_{s=1}^S \log p(\theta^{(s)}, \mathbf{x}) \in \mathbb{R}.$$

Results



Summary

- Symmetries are great learning in "deterministic models"

- Symmetries are great learning in "deterministic models"
- Symmetries are very problematic for statistical models

- Symmetries are great learning in "deterministic models"
- Symmetries are very problematic for statistical models
- "Under parametrised" approximate posteriors leads to pathological measures

- The Laplace approximation severely underfits because it does not reflect the re-parametrisations of functions

- The Laplace approximation severely underfits because it does not reflect the re-parametrisations of functions
- The Linearised Laplace approximation is infinitesimally invariant to re-parametrisations

- The Laplace approximation severely underfits because it does not reflect the re-parametrisations of functions
- The Linearised Laplace approximation is infinitesimally invariant to re-parametrisations
- The covariance of the Linearised Laplace Approximation defines a Riemannian metric on the Manifold of effective parameters

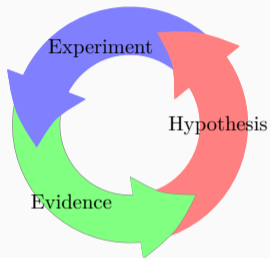
- The factorisation tells what degrees of freedom are connected to data and to the prior

- The factorisation tells what degrees of freedom are connected to data and to the prior
- What parametrisations are useful for algorithms?

- The factorisation tells what degrees of freedom are connected to data and to the prior
- What parametrisations are useful for algorithms?
- We can formulate approximate distributions that reflect the geometry of the parametrisation

- The factorisation tells what degrees of freedom are connected to data and to the prior
- What parametrisations are useful for algorithms?
- We can formulate approximate distributions that reflect the geometry of the parametrisation
- Factorisation of measures



- The factorisation tells what degrees of freedom are connected to data and to the prior
- What parametrisations are useful for algorithms?
- We can formulate approximate distributions that reflect the geometry of the parametrisation
- Factorisation of measures
- Matrix free algebra allows us to approximate parameter spaces with millions of parameters




$$p(w) = \mathcal{N}(0, \alpha \mathbf{I})$$

eof

References

-  Immer, Alexander, Maciej Korzepa, and Matthias Bauer (2021). “Improving predictions of Bayesian neural nets via local linearization”. In: *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 703–711.
-  Jacot, Arthur, Franck Gabriel, and Clément Hongler (2018). “Neural Tangent Kernel: Convergence and Generalization in Neural Networks”. In: *CoRR*.
-  MacKay, D. J. C. (1991). “Bayesian Methods for Adaptive Models”. PhD thesis. California Institute of Technology.

-  Ustyuzhaninov, Ivan, Ieva Kazlauskaite, Markus Kaiser, Erik Bodin, Neill D. F. Campbell, and Carl Henrik Ek (2020). **“Compositional uncertainty in deep Gaussian processes”**. In: *Proceedings of the Thirty-Sixth Conference on Uncertainty in Artificial Intelligence, UAI 2020, virtual online, August 3-6, 2020*. Ed. by Ryan P. Adams and Vibhav Gogate. Vol. 124. Proceedings of Machine Learning Research. AUAI Press, pp. 480–489.