

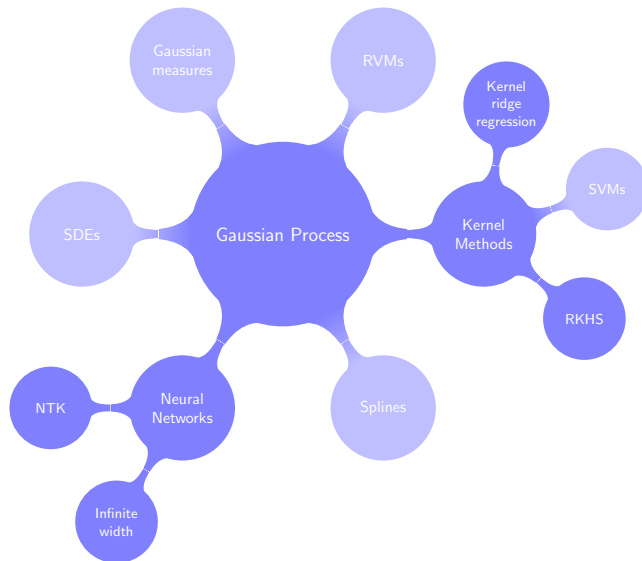
# CS-E4895 Gaussian Processes

## Lecture 8: Theory & Advanced Topics

Martin Trapp

Aalto University

Tuesday 21.03.2023



# Hilbert Space

A vector space  $\mathcal{V}$  is a set of vectors that is closed under addition and scalar multiplication.

If  $\mathcal{V}$  is equipped with a norm  $\|\cdot\|_{\mathcal{V}} \in \mathbb{R}$ , it is a *normed (vector) space*.

A Hilbert space  $\mathcal{H}$  is a complete<sup>1</sup> inner product space, with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  and induced norm  $\|x\|_{\mathcal{H}} = \sqrt{\langle x, x \rangle_{\mathcal{H}}}$ .

Recall: Operations on inner products (scalar products):

- $\langle x, x \rangle \geq 0$  and  $\langle x, \mathbf{0} \rangle = 0$
- $\langle x, y \rangle = \langle y, x \rangle$
- $\langle x + y, z \rangle = \langle x, z \rangle + \langle y, z \rangle$
- $\langle \alpha x, y \rangle = \alpha \langle x, y \rangle$

---

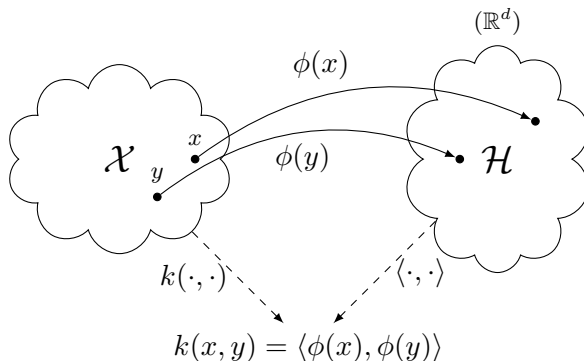
<sup>1</sup>All Cauchy sequences are convergent.

## Recap: Kernel function

Recap from the previous lectures.

A function  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a *kernel* function if and only if there exists a Hilbert space  $\mathcal{H}$  and a map  $\phi: \mathcal{X} \rightarrow \mathcal{H}$  such that:

$$k(x, y) = \langle \phi(x), \phi(y) \rangle \quad (1)$$



# Reproducing Kernel Hilbert Space

We said: Given a space  $\mathcal{X}$  and a kernel  $k$  on  $\mathcal{X}$ , *there exists* a Hilbert space  $\mathcal{H}$  and a map  $\phi$ , such that:

$$k(x, y) = \langle \phi(x), \phi(y) \rangle \quad (2)$$

for all  $x, y \in \mathcal{X}$ .

First: Let  $\phi: \mathcal{X} \rightarrow \mathbb{R}^{\mathcal{X}}$  and let us define:

$$k_x := \phi(x) = k(x, \cdot) \quad (3)$$

Therefore, we have that:  $k_x(y) = k(x, y)$ .

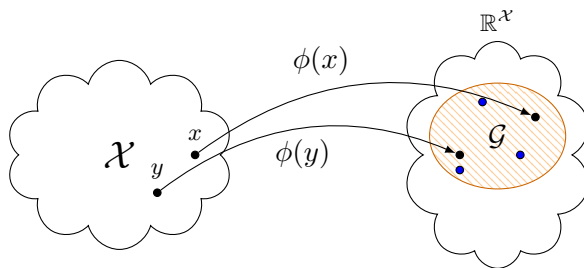
Second: Let  $\mathcal{G}$  denote a vector space with span based on the images  $\{k_x \mid x \in \mathcal{X}\}$ , i.e.,

$$\mathcal{G} := \left\{ \sum_{i=1}^m \alpha_i k_{x_i} \mid \alpha_i \in \mathbb{R}, m \in \mathbb{N}, x_i \in \mathcal{X} \right\}. \quad (4)$$

# Reproducing Kernel Hilbert Space

Let  $\mathcal{G}$  denote a vector space with span based on the images  $\{k_x \mid x \in \mathcal{X}\}$ , i.e.,

$$\mathcal{G} := \left\{ \sum_{i=1}^m \alpha_i k_{x_i} \mid \alpha_i \in \mathbb{R}, m \in \mathbb{N}, x_i \in \mathcal{X} \right\}. \quad (5)$$



# Reproducing Kernel Hilbert Space

Now: Let's define an inner product on  $\mathcal{G}$ .

Note: By definition of a kernel function, we can only choose:

$$\langle k_x, k_y \rangle := k(x, y) \quad (6)$$

(Recall  $k_x = k(x, \cdot)$ , hence,  $\langle k_x, k_y \rangle = \langle k(x, \cdot), k(y, \cdot) \rangle$ .)

Therefore, for any  $f, g \in \mathcal{G}$ , with  $f = \sum_i \alpha_i k_{x_i}$  and  $g = \sum_j \beta_j k_{y_j}$ , we have:

$$\langle f, g \rangle = \left\langle \sum_i \alpha_i k(x_i, \cdot), \sum_j \beta_j k(y_j, \cdot) \right\rangle \quad (7)$$

$$= \sum_{i,j} \alpha_i \beta_j \underbrace{\langle k_{x_i}, k_{y_j} \rangle}_{=k(x_i, y_j)} \quad (8)$$

Side note: To make  $\mathcal{G}$  a Hilbert space, we need to make it complete, *i.e.*, ensure all Cauchy sequences converge = "ensure that no points are missing".

# Reproducing Kernel Hilbert Space

## Definition (Reproducing kernel Hilbert space (RKHS))

Let  $\mathcal{H}$  be a Hilbert space of real functions  $f$  defined on an index set  $\mathcal{X}$ . Then  $\mathcal{H}$  is called a reproducing kernel Hilbert space endowed with an inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  if there exists a kernel function  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  with the following properties:

- for every  $x \in \mathcal{X}$ ,  $k_x(y) = k(x, y)$  as a function of  $y \in \mathcal{X}$  belongs to  $\mathcal{H}$ , and
- $k$  has the reproducing property.

Reproducing property:

$$\langle k_x, f \rangle = \langle k_x, \sum_i \alpha_i k_{x_i} \rangle \quad (9)$$

$$= \sum_i \alpha_i \langle k_x, k_{x_i} \rangle = \sum_i \alpha_i k(x, x_i) = f(x) \quad (10)$$

Note: Given a kernel, there is a unique RKHS. Given an RKHS, there is a unique kernel.  
(Moore-Aronszajn theorem)



# Representer Theorem

Setting:

- We are given a kernel  $k$  and denote the corresponding RKHS as  $\mathcal{H}$ .
- We want to learn a linear function  $f(\mathbf{x})$  from a finite data set  $\{\mathbf{x}_i, y_i\}_{i=1}^n$ .

## Theorem (Representer theorem)

Consider the risk minimization problem of the form:

$$\min_{f \in \mathcal{H}} \underbrace{R_n(\mathbf{y}, \mathbf{f})}_{\text{empirical risk}} + \lambda \underbrace{\Omega(\|f\|_{\mathcal{H}})}_{\text{regularizer}} \quad (11)$$

where  $\mathbf{f} = \{f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)\}$ ,  $\mathbf{y} = \{y_1, \dots, y_n\}$ , and  $\lambda$  is a scaling parameter.

Then eq. (11) always has an optimal solution of the form:

$$\boxed{f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})} \quad (12)$$

## Representer Theorem: Example

Let's consider the following risk minimization problem:

$$\min_{f \in \mathcal{H}} \underbrace{\sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2}_{=R_n(\mathbf{y}, \mathbf{f})} + \lambda \underbrace{\|f\|_{\mathcal{H}}^2}_{=\Omega(\|f\|_{\mathcal{H}})} . \quad (13)$$

This minimization problem is known as *kernel ridge regression*.

Let's plug in the solution according to the *representer theorem* ( $f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$ ):

$$\min_{f \in \mathcal{H}} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|_{\mathcal{H}}^2 \quad (14)$$

$$= \min_{\alpha} \sum_{i=1}^n (y_i - \sum_{j=1}^n \alpha_j k(\mathbf{x}_i, \mathbf{x}_j))^2 + \lambda \left\| \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \cdot) \right\|_{\mathcal{H}}^2 \quad (15)$$

## Representer Theorem: Example

Let  $\mathbf{K}$  denote the kernel matrix  $\mathbf{K}_{\mathcal{X},\mathcal{X}}$ ,  $\mathbf{y} = (y_1, \dots, y_n)^\top$ , then we can write the objective as:

$$J(\boldsymbol{\alpha}) = \sum_{i=1}^n \left( y_i - \sum_{j=1}^n \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) \right)^2 + \lambda \left\| \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \cdot) \right\|_{\mathcal{H}}^2 \quad (16)$$

$$= \boldsymbol{\alpha}^\top (\mathbf{K}\mathbf{K}^\top + \lambda \mathbf{K}) \boldsymbol{\alpha} - 2\mathbf{y}^\top \mathbf{K} \boldsymbol{\alpha} \quad (17)$$

Note that  $J(\boldsymbol{\alpha})$  is convex and  $\mathbf{K}$  is assumed positive-definite (hence, invertible).

Then, by taking  $\frac{\partial J(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}} = 0$  we obtain:

$$\hat{\boldsymbol{\alpha}} = (\mathbf{K} + \lambda I)^{-1} \mathbf{y} \quad (18)$$

and the prediction for test point  $\mathbf{x}^*$  is, therefore,

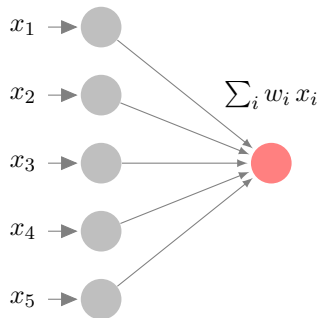
$$\hat{f}(\mathbf{x}^*) = \sum_i^n \hat{\alpha}_i k(\mathbf{x}_i, \mathbf{x}^*) = \underbrace{\mathbf{k}^\top (\mathbf{K} + \lambda I)^{-1} \mathbf{y}}_{\text{predictive mean of GP regression}} \quad (19)$$

$\implies$  Minimizer to kernel ridge regression = predictive mean of GP regression.

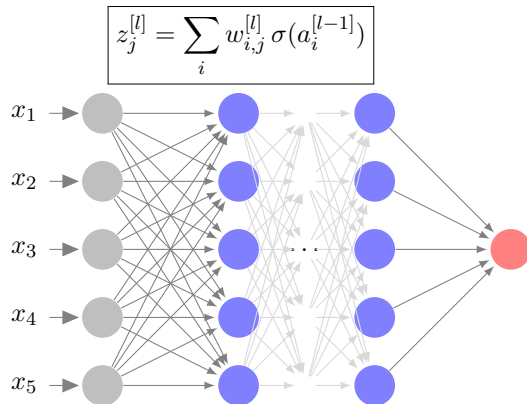
## Further readings

- C. E. Rasmussen & C. K. I. Williams, *Gaussian Processes for Machine Learning*, MIT press 2006, Chapter 6.
- A. Gretton, *Introduction to RKHS, and some simple kernel algorithms*, Lecture at UCL 2013.
- C. Heil, *Banach and hilbert space review.*, tech. rep., Georgia Tech, 2006.

Perceptron



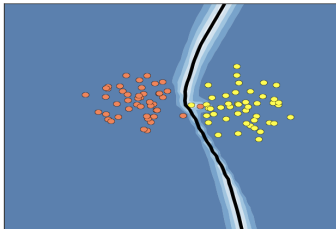
Multilayer Perceptron / NN



- $\sigma(\cdot)$  is commonly chosen to be non-linear, e.g., rectified linear unit (ReLU) defined as  $\sigma(x) = \max\{x, 0\}$ .
- The deep learning community extends this concept by a battery of modules & techniques which we will not discuss here.

# Neural Network Classification: Example

Deterministic Neural Network



Bayesian Neural Network (BNN)

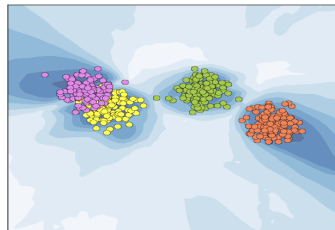
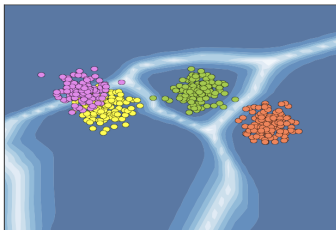
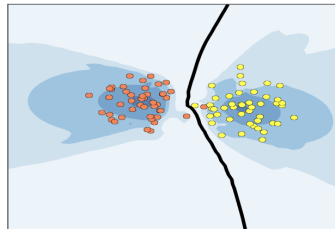
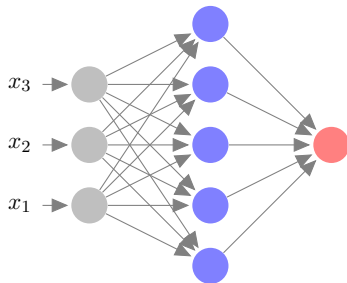


Illustration by: A. Kristiadi, M. Hein, and P. Hennig. *Being Bayesian, Even Just a Bit, Fixes Overconfidence in ReLU Networks*, ICML 2020.

# Infinitely-wide single layer BNNs

First consider a single-hidden layer neural network of the form:

$$z_i^{[1]}(\mathbf{x}) = b_i^{[1]} + \sum_{j=1}^{K_1} w_{i,j}^{[1]} a_j^{[1]}(\mathbf{x}), \quad a_j^{[1]}(\mathbf{x}) = \sigma \left( b_j^{[0]} + \sum_{d=1}^D w_{j,d}^{[0]} x_d \right) \quad (20)$$



Assume the following generative process for the parameters:

$$w_{i,j}^{[1]} \sim \mathcal{N}(0, \frac{\sigma_w^2}{K_1}), \quad b_i^{[1]} \sim \mathcal{N}(0, \sigma_b^2), \quad w_{j,d}^{[0]} \sim \mathcal{N}(0, \frac{\sigma_w^2}{D}), \quad b_j^{[0]} \sim \mathcal{N}(0, \sigma_b^2) \quad (21)$$

# Infinitely-wide single layer BNNs

Since  $z_i^{[1]}(\mathbf{x})$  is a sum of RVs with finite moments, then as  $K_1 \rightarrow \infty$  the network output  $z_i^{[1]}(\mathbf{x})$  converges in distribution to:

$$b_i^{[1]} + w_{i,1}^{[1]}a_1^{[1]}(\mathbf{x}) + w_{i,2}^{[1]}a_2^{[1]}(\mathbf{x}) + \dots \xrightarrow{\mathcal{D}} \sqrt{1 + K_1}(Z_i(\mathbf{x}) - \mu) \quad (22)$$

with  $Z_i(\mathbf{x}) \sim \mathcal{N}(0, \sigma^2)$ . (CLT)

Following the multivariate CLT, we have that  $\mathbf{Z}_i = (Z_i(\mathbf{x}_1), Z_i(\mathbf{x}_2), \dots, Z_i(\mathbf{x}_n))^T$  is joint multivariate Gaussian distributed.

Recall: A Gaussian process (GP) is a collection of RVs  $\mathbf{F}$  indexed by  $\mathcal{X}$ , where any finite subset of  $\mathbf{F}$  is joint multivariate Gaussian distributed, and of which any two overlapping finite sets are marginally consistent.

$\implies$  A single-hidden layer BNN (prior) converges to a GP when  $K_1 \rightarrow \infty$ . (Neal 1994)



# Infinitely-wide single layer BNNs

First question: What is the mean of this process?

$$\mathbb{E}[Z_i(\mathbf{x})] = \mathbb{E}[b_i^{[1]} + \sum_{j=1}^{K_1} w_{i,j}^{[1]} a_j^{[1]}(\mathbf{x})] \quad (23)$$

$$= \underbrace{\mathbb{E}[b_i^{[1]}]}_{=0} + \underbrace{\mathbb{E}[\sum_{j=1}^{K_1} w_{i,j}^{[1]} a_j^{[1]}(\mathbf{x})]}_{=\mathbf{0}^\top \mathbf{a}^{[1]}(\mathbf{x})=0} \quad (24)$$

Second question: What is the induced kernel?

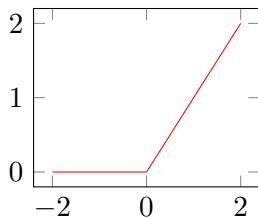
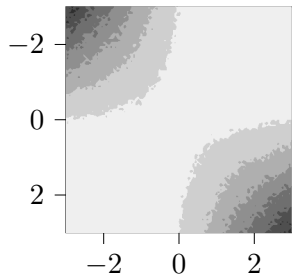
For this, let's examine the covariance:

$$\text{cov}(Z_i(\mathbf{x}_1), Z_i(\mathbf{x}_2)) = \mathbb{E}[Z_i(\mathbf{x}_1)Z_i(\mathbf{x}_2)] - \underbrace{\mathbb{E}[Z_i(\mathbf{x}_1)]\mathbb{E}[Z_i(\mathbf{x}_2)]}_{=0} \quad (25)$$

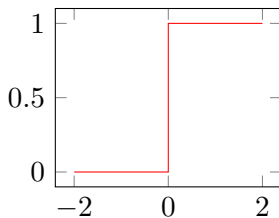
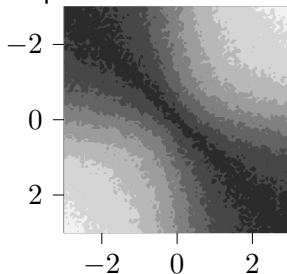
$$= \sigma_b^2 + \sigma_w^2 \underbrace{\mathbb{E}[a_i^1(\mathbf{x}_1)a_i^1(\mathbf{x}_2)]}_{=k(\mathbf{x}_1, \mathbf{x}_2)} \quad (26)$$

# Infinitely-wide single layer BNNs

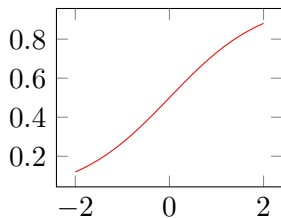
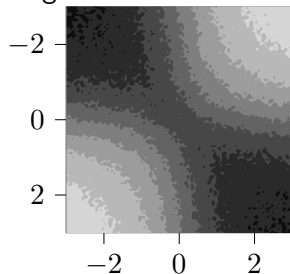
ReLU  $\sim$  ArcCos-1 kernel



Step  $\sim$  ArcCos-0 kernel



Sigmoid  $\sim$  NN kernel



# Infinitely-wide deep BNNs

Extensions to multiple hidden layers (aka deep neural networks).

$$z_i^{[1]}(\mathbf{x}) = b_i^{[1]} + \sum_{j=1}^{K_1} w_{i,j}^{[1]} a_j^{[1]}(\mathbf{x}), \quad a_j^{[1]}(\mathbf{x}) = \sigma \left( b_j^{[0]} + \sum_{d=1}^D w_{j,d}^{[0]} x_d \right) \quad (27)$$

Consider  $L$  layers with *i.i.d.* parameters, then  $z_j^{[l-1]}(\mathbf{x})$  are *i.i.d.* draws from a GP.

Therefore,  $z_j^{[l]}(\mathbf{x})$  is a sum of *i.i.d.* terms and

$$b_i^{[l]} + w_{i,1}^{[l]} a_1^{[l]}(\mathbf{x}) + w_{i,2}^{[l]} a_2^{[l]}(\mathbf{x}) + \dots \xrightarrow{\mathcal{D}} \sqrt{1 + K_l} (Z_i^{[l]}(\mathbf{x}) - \mu) \quad (28)$$

with  $Z_i^{[l]}(\mathbf{x}) \sim \mathcal{N}(0, \sigma^2)$  (CLT) and  $\mathbf{Z}_i^{[l]} = (Z_i^{[l]}(\mathbf{x}_1), Z_i^{[l]}(\mathbf{x}_2), \dots, Z_i^{[l]}(\mathbf{x}_n))^{\top}$  is joint multivariate Gaussian distributed (multivariate CLT).

$\implies$  Covariance structure now recursively defined.

# Infinitely-wide deep BNNs

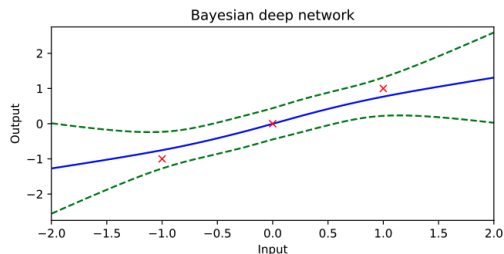
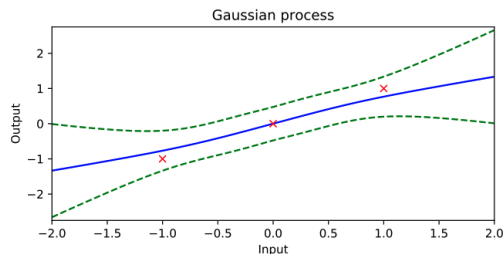


Illustration by: A. G. Matthews et al. *Gaussian process behaviour in wide deep neural networks*, ICLR 2018.

Note: Kernels obtained from infinite-width BNNs are fixed (not trainable), and do not accurately reflect the representation learning of NNs.

## Further reading

- R. Neal *Bayesian Learning for Neural Networks*, PhD thesis, 1994.
- J. Lee et al. *Deep neural networks as Gaussian processes*, ICLR 2018.
- A. G. Matthews et al. *Gaussian process behaviour in wide deep neural networks*, ICLR 2018.

# Neural Tangent Kernel

The Neural Tangent Kernel (NTK) establishes a link between gradient descent training in NNs and kernel methods.

Setting:

$$f(\mathbf{x}; \theta) = \frac{1}{\sqrt{K}} \sum_{j=1}^K w_j^{[1]} \sigma \left( \sum_{d=1}^D w_{j,d}^{[0]} x_d \right) \quad (29)$$

where  $\theta$  denotes all parameters (weights and biases).

Objective:

$$R_n = \frac{1}{2} \sum_{i=1}^n (f(\mathbf{x}_i; \theta) - y_i)^2 \quad (30)$$

$$R_n = \frac{1}{2} \sum_{i=1}^n (f(\mathbf{x}_i; \theta) - y_i)^2 \quad (31)$$

Gradient Descent (GD):

$$\theta_{t+1} = \theta_t - \eta \sum_{i=1}^n (f(\mathbf{x}_i; \theta) - y_i) \nabla_{\theta} f(\mathbf{x}_i; \theta_t) \quad (32)$$

where  $\eta$  is a learning rate that is usually small.

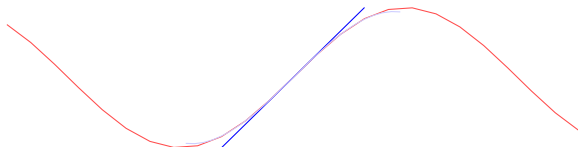
- If  $f(\mathbf{x}_i; \theta_t)$  is linear:  $\Rightarrow \nabla_{\theta} f(\mathbf{x}_i; \theta_t)$  is constant (only depends on  $\mathbf{x}_i$ )
- If  $f(\mathbf{x}_i; \theta_t)$  is non-linear:  $\Rightarrow \nabla_{\theta} f(\mathbf{x}_i; \theta_t)$  is changing over time

Empirical observation: If  $K$  is large then the parameters  $\theta$  of a NN learned with GD do not change much over time. (lazy training)

# Neural Tangent Kernel

Lazy training motivates: First-order Taylor approximation of  $f(\mathbf{x}, \theta)$  around  $\theta_0$  (initialization):

$$f(\mathbf{x}; \theta) \approx f(\mathbf{x}; \theta_0) + \nabla_{\theta} f(\mathbf{x}; \theta_0)^{\top} (\theta - \theta_0) + \dots \quad (33)$$



Note:  $\nabla_{\theta} f(\mathbf{x}; \theta_0)$  does not depend on  $\theta$  and is nonlinear in  $\mathbf{x}$ .

Trick: We define  $\phi(\mathbf{x}) := \nabla_{\theta} f(\mathbf{x}; \theta_0)$ .

$$k(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = \langle \nabla_{\theta} f(\mathbf{x}_i; \theta_0), \nabla_{\theta} f(\mathbf{x}_j; \theta_0) \rangle \quad (34)$$



# Neural Tangent Kernel

The NTK allows us to derive the induced kernel of a infinite-width NN (often in closed-form).

We can also analyse the training dynamics:

$$\frac{\theta_{t-1} - \theta_t}{\eta} = -\nabla_{\theta} R_n(\theta_t), \quad R_n(\theta) = \frac{1}{2} \sum_{i=1}^n (f(\mathbf{x}_i; \theta) - y_i)^2 \quad (35)$$

As we have  $\eta \rightarrow \infty$ : (gradient flow)

$$\frac{d\theta(t)}{dt} = -\nabla_{\theta} R_n(\theta(t)) = -\nabla_{\theta} \hat{\mathbf{y}}(\theta(t)) (\hat{\mathbf{y}}(\theta(t)) - \mathbf{y}) \quad (36)$$

$$\frac{d\hat{\mathbf{y}}(\theta(t))}{dt} = - \underbrace{\nabla_{\theta} \hat{\mathbf{y}}(\theta(t))^{\top} \nabla_{\theta} \hat{\mathbf{y}}(\theta(t))}_{=\mathbf{K}_{\text{NTK}}} (\hat{\mathbf{y}}(\theta(t)) - \mathbf{y}) \quad (37)$$

## Further readings

- A. Jacot, F. Gabriel, and C. Hongler, *Neural tangent kernel: Convergence and generalization in neural networks*, NeurIPS 2018.
- L. Chizat, F. Bach, *A Note on Lazy Training in Supervised Differentiable Programming*, tech. rep., INRIA 2019.
- <https://github.com/kwignb/NeuralTangentKernel-Papers>

# Summary

Recap:

- Theory on RKHS allows us to better understand kernel functions.
- The representer theorem  $\Rightarrow$  Optimal solution:  $\hat{f}(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$
- Kernel ridge regression = GP predictive mean
- BNN (at initialization) converges to a GP at infinite-width limit
- NTK establishes a link between NN training with GD and kernel methods

Advances in Probabilistic Machine Learning

<https://aaltoml.github.io/apml/>

CS-E407516: Special Course in ML, DS & AI: **Tractable Probabilistic Modelling**

<https://mycourses.aalto.fi/course/view.php?id=38446>