


STAT 535

12/8/22

Lecture 20

(... and last)



- Wide Neural networks as GP
- The NTK
- Deep NN \xrightarrow{R} Class

Gaussian Process

$\mathcal{X} \ni x$ distribution over functions on \mathcal{X}

- $f(x) \sim N(0, k(x, x))$ for all x

↑
or known

- any set $x^{1:n}$, $f(x^{1:n}) \sim N(0, G_x)$

define $k(x, x') = \text{cov}(f(x), f(x')) \equiv \underline{E[f(x)f(x')]}$

↑
Mercer kernel

Assume $k(x, x') = K(\|x - x'\|)$ e.g. Gaussian kernel

• Pb

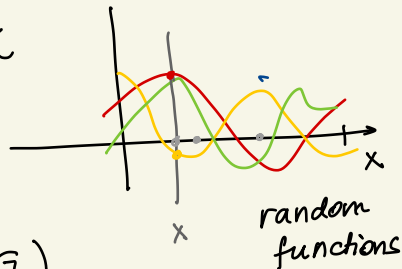
know K

observe $(x^{1:n}, y^{1:n} = f(x^{1:n})) = \tilde{d}$

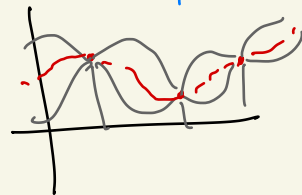
Bayesian: Prior $f \sim GP(0, K)$

Posterior $f | \tilde{d} = ?$

Prior $GP(0, K)$



Posterior
 $GP(\mu_{x^{1:n}}, K_{x^{1:n}})$



$$\Sigma \equiv \kappa_{xx,xx} = \begin{bmatrix} G & \cdot \\ \cdot & \cdot \end{bmatrix} \begin{matrix} \left. \vphantom{\begin{bmatrix} G & \cdot \\ \cdot & \cdot \end{bmatrix}} \right\} x^{1:n} \\ \left. \vphantom{\begin{bmatrix} G & \cdot \\ \cdot & \cdot \end{bmatrix}} \right\} x \\ \kappa(x,x) \end{matrix}$$

Same problem:

→ some $x \in \mathcal{X}$: wanted $f(x|\mathcal{D}) \sim N(\mu, \sigma^2)$
posterior

\uparrow \uparrow
 $\mu_{x|\mathcal{D}}$ $\sigma^2_{x|\mathcal{D}}$

Remark $\begin{bmatrix} y \\ y^{1:n} \end{bmatrix} \sim N(0_{n+1}; \Sigma)$ } → want conditional
jointly Gaussian $y | y^{1:n}, x^{1:n} \sim N(\mu_{x|\mathcal{D}}, \sigma^2_{x|\mathcal{D}})$

$$\rightarrow \mu_{x|\mathcal{D}} = \begin{matrix} \xrightarrow{\kappa(x, X)} \\ \underbrace{[\kappa(x, x') \dots]}_n \end{matrix} \begin{matrix} \underbrace{G^{-1}}_{n \times n} \end{matrix} \begin{matrix} \underbrace{[y^{1:n}]}_n \end{matrix}$$

$$\sigma^2_{x|\mathcal{D}} = \underbrace{\kappa(x, x)}_{\text{red}} - \underbrace{\kappa(x, X)}_{\text{orange}} \underbrace{G^{-1}}_{n \times n} \underbrace{\kappa(X, x)}_{\text{yellow}}$$

Lecture VII – Wide multilayer networks and the Neural Tangent Kernel (NTK)

Marina Meilă
`mmp@stat.washington.edu`

Department of Statistics
University of Washington

October, 2022

The Neural Tangent Kernel (NTK)

Wide networks and Gaussian Processes

The NTK is constant during training

Example – regression and \mathcal{L}_{LS}

Wide and deep networks and classification

Notation

- ▶ Neural network predictor $f(x; \theta)$, where $x \in \mathbb{R}^d$
- ▶ For each layer $l = 1 : L$ of dimension m_l , with $x^0 \equiv x$, and $z^L \equiv f(x)$

$$z^{l+1} = W^{l+1}x^l + b^{l+1} \quad x^{l+1} = \phi(z^{l+1}) \quad (1)$$

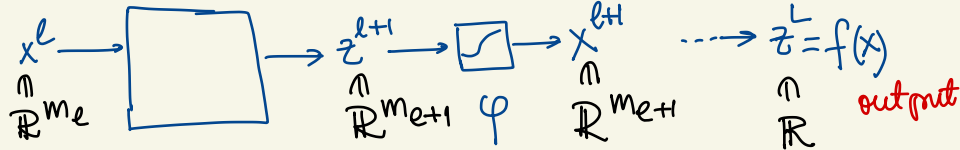
Here $x^{l,l+1}, z^{l+1}, b^{l+1}$ are column vectors W^{l+1} is a $m_{l+1} \times m_l$ matrix, $\phi()$ is the non-linearity/activation function.

- ▶ The weights

$$W_{ij}^l = \sigma_w w_{ij}^l / \sqrt{m_l}, \quad b_j^l = \sigma_b \beta_j^l, \quad \text{Known as NTK parametrization} \quad (2)$$

- ▶ Parameter vector $\theta = \text{vector}\{w^{1:L}, \beta^{1:L}\} \in \mathbb{R}^p$ initialized i.i.d. $\sim N(0, 1)$
- ▶ $\sigma_{w,b}$ are fixed hyper-parameters, $1/\sqrt{m_l}$ normalizes the expected norm of W^l columns
- ▶ Loss $\mathcal{L}(y, f)$
- ▶ We want to analyze the behavior of this network $f()$ at initialization and during training, when $m_{1:L}$ very large
- ▶ Three approximations help analysis
 - (A1) continuous time training, called **gradient flow**
 - (A2) $m_{1:L} \rightarrow \infty$ in the wide limit, we can apply the Central Limit Theorem (CLT), and Gaussian Processes
 - (A3) parameters θ do not change much during training, i.e. $\theta_t - \theta_0$ is small

$x^0 = x$
input



$$z^{l+1}(x) = W^{l+1} x^l + b^{l+1}$$

$$l = 1:L$$

$$z_j^{l+1}$$

$$j = 1:m_{l+1}$$

Initialization

$$W_{jv}^{l+1} \sim N(0, \frac{1}{m_l} \underbrace{\sigma_w^2}_{\text{variance}})$$

$$b_j^{l+1} \sim N(0, \underbrace{\sigma_b^2}_1)$$

1

1) x fixed, θ random
 $\theta = (W^{1:L}, b^{1:L})$

$$z_j^{l+1}(x) \xrightarrow{\text{CLT}} z^{l+1} \sim N(0, \quad)$$

0) $x_j^l(x)$ k.v. $\in \mathbb{R}^{m_l}$ iid

$m_{1:L} \rightarrow \infty$ wide network

2) $z^l \sim N(0, \quad)$ multivariate gaussian

No proof 3) $z^l(x) \sim \text{GP on } \mathcal{X}$

The NTK

Neural Tangent Kernel

K^L for GPNN

layer l : $z^l \sim \text{GP}(0, K^l)$ (at init for now)

$f \equiv z^L \sim \text{GP}(0, K^L)$

training

The Gradient Flow

- Assume training by **gradient descent** on $\hat{\mathcal{L}} = \sum_i \mathcal{L}(y^i, f(x^i))$
- The gradient of $\hat{\mathcal{L}}$

empirical loss

$$\mathcal{L}_{\mathcal{D}} = \left[\mathcal{L}(y^i, f(x^i)) \right]_{i=1:n} \in \mathbb{R}^n$$

$$\nabla_{\theta} \hat{\mathcal{L}} = \sum_i \underbrace{\frac{\partial \mathcal{L}}{\partial f}}_{\text{green}}(y^i, f(x^i; \theta)) \underbrace{\nabla_{\theta} f(x^i, \theta)}_{\text{red}} = \underbrace{\nabla_{\theta} f_{\mathcal{D}}}_{\text{red}} \underbrace{\nabla_f \mathcal{L}_{\mathcal{D}}}_{\text{green}} \in \mathbb{R}^p \quad (3)$$

where $\nabla_f \mathcal{L}_{\mathcal{D}} = \left[\frac{\partial \mathcal{L}}{\partial f}(y^i, f(x^i; \theta)) \right]_{i=1:n} \in \mathbb{R}^n$, $\nabla_{\theta} f_{\mathcal{D}} = [\nabla_{\theta} f(x^i, \theta)]_{i=1:n} \in \mathbb{R}^{p \times n}$

- Assume (A1) gradient descent with infinitesimal time steps. In other words, the parameters evolve by an ordinary differential equation

$\frac{d\theta}{dt}$:

$$\dot{\theta} = -\eta \underbrace{\nabla_{\theta} f_{\mathcal{D}}}_{\text{red}} \underbrace{\nabla_f \mathcal{L}_{\mathcal{D}}}_{\text{green}} \in \mathbb{R}^p$$

$$\dot{f}(x) = \sum_{j=1}^p \frac{\partial f}{\partial \theta_j} \frac{\partial \theta_j}{\partial t} = \underbrace{(\nabla_{\theta} f)}_{\text{red}}^T \underbrace{\dot{\theta}}_{\text{blue}} \in \mathbb{R} \quad (4)$$

$$f_{\mathcal{D}} = \left[f(x^i) \right]_{i=1:n}$$

$$\dot{f}_{\mathcal{D}} = -\eta \underbrace{(\nabla_{\theta} f_{\mathcal{D}})^T}_{\text{red}} \underbrace{\nabla_{\theta} f_{\mathcal{D}}}_{\text{red}} \underbrace{\nabla_f \mathcal{L}_{\mathcal{D}}}_{\text{green}} \in \mathbb{R}^p \quad (5)$$

$$\nabla_{\theta_j} f(x^i)$$

- $G \equiv \nabla_{\theta} f_{\mathcal{D}}^T \nabla_{\theta} f_{\mathcal{D}} \equiv \kappa(X, X)$ is a **Gram matrix**!
- Therefore, we define the **Neural Tangent Kernel (NTK)** by

Gram matrix

$$\kappa(x, x') = \nabla_{\theta} f(x; \theta)^T \nabla_{\theta} f(x'; \theta) \quad (7)$$

$$G_{i,j} = \left[\kappa(x^i, x^j) \right]_{i,j=1:n}$$

tangent to f in θ space

Gradient flow and NTK – summary

$$\begin{aligned}\dot{\theta} &= -\eta \nabla_{\theta} f_{\mathcal{D}} \nabla_f \mathcal{L}_{\mathcal{D}} \in \mathbb{R}^p \\ \dot{f}_{\mathcal{D}} &= -\eta G \nabla_f \mathcal{L} \in \mathbb{R}^p \\ \kappa(x, x') &= \nabla_{\theta} f(x)^T \nabla_{\theta} f(x')\end{aligned}$$

- ▶ $f_X, \nabla_{\theta} f_X, G$ depend only on the inputs X, θ
- ▶ $\nabla_f \mathcal{L}_{\mathcal{D}}$ depends only on the correct outputs Y , and predicted outputs, i.e. on Y and θ
- ▶ This holds for **any predictor**! So what is special about neural networks?

Gradient flow and NTK – summary

$$\begin{aligned}\dot{\theta} &= -\eta \nabla_{\theta} f_{\mathcal{D}} \nabla_{\theta} \mathcal{L}_{\mathcal{D}} && \in \mathbb{R}^p \\ \dot{f}_{\mathcal{D}} &= -\eta G \nabla_{\theta} \mathcal{L} && \in \mathbb{R}^p \\ \kappa(x, x') &= \nabla_{\theta} f(x)^T \nabla_{\theta} f(x')\end{aligned}$$

- ▶ $f_X, \nabla_{\theta} f_X, G$ depend only on the inputs X, θ
- ▶ $\nabla_{\theta} \mathcal{L}$ depends only on the correct outputs Y , and predicted outputs, i.e. on Y and θ
- ▶ This holds for **any predictor!** So what is special about neural networks?
- ▶ First, we will analyze κ for very wide neural networks with random parameters (e.g. at initialization)
- ▶ Then, we will analyze what happens during training under assumption (A3)

iii

i $\equiv \kappa^L$ of GP
ii $m_{1:L} \rightarrow \infty \Rightarrow \kappa^{1:L}$ deterministic

Wide NN's Gaussian Process (GP)

- ▶ This is about f_0 , a NN initialized with Gaussian independent parameters. For simplicity, we denote it as f .
- ▶ Assume $\theta^{1:L-1}$ fixed, only W^L, b^L random as in (??)
- ▶ Recall $f(x) = W^L x^{L-1}(x) + b^L$ for any x with $x^{L-1} \in \mathbb{R}^{m_L}$
- ▶ $f(x)$ = sum of m_{L-1} i.i.d. random variables, hence $f(x) \sim \text{Normal}$ by CLT, for m_{L-1} large
- ▶ Randomness is over weights W^L, b^L !!!
- ▶ We have $E[f(x)] = 0$ and

$$\text{Cov}(f(x), f(x')) = E[(W^L x^{L-1} + b^L)(W^L (x')^{L-1} + b^L)] = \frac{\sigma_w^2}{m_{L-1}} (x^{L-1})^T (x')^{L-1} + \sigma_b^2 \equiv \kappa^L(x, x') \quad (8)$$

where $x^{L-1}, (x')^{L-1} \in \mathbb{R}^{m_{L-1}}$ are the outputs of the $(L-1)$ -th layer for inputs x, x'

- ▶ κ^L is a positive definite **kernel** Exercise Prove this.
- ▶ $f(x)$ is a random function of x
- ▶ The distribution of $f(x)$ defined as above, is called a **Gaussian Process**
- ▶ More generally, it can be shown [Jacot, Gabriel, Hongler, NeurIPS 2018] that, when all θ parameters are sampled as in (??), $f_0(x) \sim GP(0, \kappa^L)$

Q1 What is the kernel κ^L of this GP

Q2 This is all nice, but θ changes during training. What can we say about θ_t, f_t after training?

Q1: Idea.

i

- From (8), for layer $l = 1 : L$ we have

$$\kappa^l(x, x') = E[z_j^l(x) z_j^l(x')] = \frac{\sigma_w^2}{m_{l-1}} (x^{l-1})^T (x')^{l-1} + \sigma_b^2 \quad (9)$$

with $x^{l-1} = \phi(z^{l-1})$. Note also that z_j^l are i.i.d. so it does not matter which j we choose.

- In particular, $\kappa^1(x, x') = \frac{\sigma_w^2}{m_1} x^T x' + \sigma_b^2$ is deterministic
- ... and κ^l is random for $l > 1$.
- However, when $m_l \rightarrow \infty$, $\frac{1}{m_{l-1}} (x^{l-1})^T (x')^{l-1} \rightarrow E[*]$
- More specifically, this expectation can be written as

$$E[*] = \int \int \phi(z) \phi(z') \text{Normal}\left(\begin{bmatrix} z \\ z' \end{bmatrix}; 0, \kappa_{x, x'}^{l-1}\right) dz dz'. \quad (10)$$

In the above z, z' represent the $z^{l-1}(x), z^{l-1}(x')$ variables, sampled from the level l Normal distribution, which has covariance given by κ^{l-1} , namely

$$\kappa_{x, x'}^{l-1} = \begin{bmatrix} \kappa^{l-1}(x, x) & \kappa^{l-1}(x, x') \\ \kappa^{l-1}(x', x) & \kappa^{l-1}(x', x') \end{bmatrix}. \quad (11)$$

- Hence, the limit of $\kappa^l(x, x')$ when $m_{1:l} \rightarrow \infty$, is a **deterministic kernel** for all l . [Jacot, Gabriel, Hongler, NeurIPS 2018] derived this recursion (next page).

Q1: A recursive expression for the Neural Tangent Kernel

[Jacot, Gabriel, Hongler, NeurIPS 2018]

- ▶ L fixed, $m \rightarrow \infty$
- ▶ Simplified expression for $m_{0:L} = m$, $\sigma_w = \sigma_b = 1$
- ▶ Then the NTK $\kappa \equiv \kappa^L$ is defined recursively by layer

$$\kappa^1(x, x') = \Sigma^1(x, x'), \quad \Sigma^1(x, x') = \frac{1}{m} x^T x' + 1 \quad (12)$$

$$\kappa^{l+1}(x, x') = \kappa^l(x, x') \dot{\Sigma}^{l+1}(x, x') + \Sigma^{l+1}(x, x'), \quad (13)$$

$$\text{with} \quad \Sigma^{l+1}(x, x') = L_{\Sigma^l(x, x')}^{\phi}, \quad (14)$$

$$\dot{\Sigma}^{l+1}(x, x') = L_{\Sigma^l(x, x')}^{\phi'}, \quad (15)$$

$$\text{and} \quad L_{\Sigma}^{\phi} = E[\phi(X)\phi(X')] \text{ with } (X, X') \sim N(0, \begin{bmatrix} \Sigma(X, X) & \Sigma(X, X') \\ \Sigma(X, X') & \Sigma(X', X') \end{bmatrix}) \quad (16)$$

- ▶ In other words, at level $l+1$, $X \equiv x^l, X' \equiv (x')^l$ are sampled from a GP with kernel Σ^l , and $\Sigma^{l+1}(x, x')$, $\dot{\Sigma}^{l+1}(x, x')$ represent their (scalar) covariance after passing through the non-linearities ϕ, ϕ' (where ϕ' is the derivative of ϕ)

Summary so far

- Now, we understand the random initialization of wide networks, with L layers.

$$f_0 \sim GP(0, \kappa^L) \quad (17)$$

where κ^L is a kernel that depends only on ϕ (and $\sigma_{b,w}^2$)

What next?

- Analysis of training by linearization iii
- Then, the NTK limit for $L \rightarrow \infty$ and its relevance for classification and regression

iv

The Linearized Network f^{lin}

after training
 $\|\theta_t - \theta_0\|$ small

Notation: $\theta_{0,t}, f_{0,t}$ = parameters, predictor at times 0, t

- Here we use (A3), the assumption that the parameters θ change little during training. Extensive evidence supports this assumption.
- First order Taylor expansion of f_t around f_0

$$f_t^{\text{lin}}(x) = f_0(x) + \nabla_{\theta} f_0(x)^T (\theta_t - \theta_0) \approx f^{\text{trained}} \quad (18)$$

non-linear in x , linear in θ

$$\nabla_{\theta} f_t^{\text{lin}} = \nabla_{\theta} f_0 \quad (19)$$

$$\kappa(x, x') = \nabla_{\theta} f_0(x)^T \nabla_{\theta} f_0(x') \quad \text{constant during training} \quad (20)$$

$$G_0 \equiv \kappa_{X,X} = \text{NTK of random net} \quad (21)$$

$$\dot{\theta}_t = -\eta \nabla_{\theta} f_0(X)^T \nabla_f \mathcal{L}(Y, f_t^{\text{lin}}(x)) \quad (22)$$

$$\dot{f}_t^{\text{lin}}(x) = -\eta \underbrace{\kappa(x, X) G_0^T}_{\text{depends on } \theta_0} \underbrace{\nabla_f \mathcal{L}(Y, f_t^{\text{lin}}(x))}_{\text{depends on } \theta_t} \quad (23)$$

NTK during training – empirical evidence

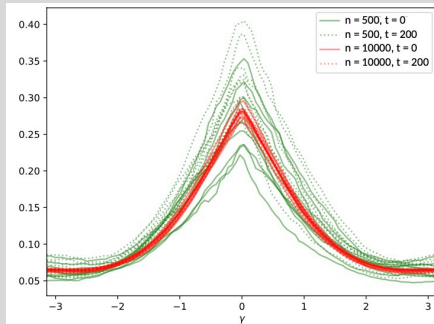


Figure 1: Convergence of the NTK to a fixed limit for two widths n and two times t .

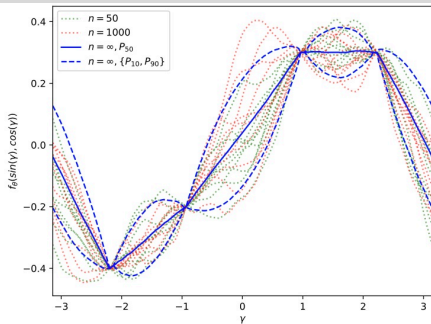


Figure 2: Networks function f_{θ} near convergence for two widths n and 10th, 50th and 90th percentiles of the asymptotic Gaussian distribution.

Linearized Network dynamics for \mathcal{L}_{LS}

- For example, for $\mathcal{L}_{LS}(y, f) = \frac{1}{2}(f - y)^2$, $\nabla_f \mathcal{L}_{LS}(f, y) = f - y$. In this case, equations (??) are a linear system and have an analytic solution.

$$\theta_t - \theta_0 = -\nabla_{\theta} f_0(X)^T G_0^{-1} \left(I - e^{-\eta G_0 t} \right) (f_0(X) - Y) \quad (24)$$

$$f_t^{\text{lin}}(X) = \left(I - e^{-\eta G_0 t} \right) Y + e^{-\eta G_0 t} f_0(X) \quad (25)$$

$$f_t^{\text{lin}}(x) = \underbrace{\kappa(x, X)^T G_0^{-1} \left(I - e^{-\eta G_0 t} \right) Y}_{\mu(x)} + \underbrace{f_0(x) - \kappa(x, X)^T G_0^{-1} \left(I - e^{-\eta G_0 t} \right) f_0(X)}_{\gamma(x)} \quad (26)$$

Notes:

- if $G_0 \succ 0$ then $e^{-\eta G_0 t} \rightarrow 0$ for $t \rightarrow \infty$
- in discrete time $t = 0, 1, 3, \dots$ replace e^{at} with $(1 - a)^t$.
Sketch of proof: $\ln(1 - a)^t = t \ln(1 - a) \approx t(-a)$ for a small; therefore $e^{-at} \approx (1 - a)^t$.
- $f_t^{\text{lin}}(x) = f_0(x) + \kappa(x, X)^T G_0^{-1} \left(I - e^{-\eta G_0 t} \right) (Y - f_0(X))$

Exercise Prove (??) from (??)

Wide and deep neural networks for classification – Basic quantities and assumptions

[Radhakrishnan, Belkin, Ullher, 2022]

- ▶ This paper studies the limits of wide neural networks $m_l \rightarrow \infty$ for all $l = 1 : L$ when the depth $L \rightarrow \infty$
- ▶ It is already known that for regression $L \rightarrow \infty$ is NOT OPTIMAL
- ▶ Since the NTK depends only of the activation function ϕ , the limit shall only depend on ϕ as well.
- ▶ In particular, the limit depends on ϕ only through the following

$$\begin{aligned} A &= E[\phi(Z)] && \text{when } z \sim N(0, 1) \\ A' &= E[\phi'(Z)] && \text{when } z \sim N(0, 1) \\ B &= E[(\phi'(Z))^2] && \text{when } z \sim N(0, 1) \end{aligned}$$

- ▶ Classifier $f(x) = \lim_{L \rightarrow \infty} \text{sgn } Y G^{-1} \kappa^L(X, x)$ with $G = [\kappa^L(x^i, x^j)]_{i,j=1:n}$.
- ▶ Additional assumptions
 - ▶ Data $X \subseteq \mathcal{S}_+^d$, vectors of norm 1 with all entries ≥ 0 .
 - ▶ Simplifying assumptions on NTK parameters (e.g. $\sigma_w = \sigma_b = 1$)

" $\kappa^L \rightarrow 0$ " $\Leftarrow L \rightarrow \infty$
↑
informal statement

Case $A \neq 0$: Networks implement majority vote

Theorem (Proposition 1 in [Radhakrishnan, Belkin, Ulher, 2022])

If there is a function $0 < c(L) < \infty$ so that

$$\lim_{L \rightarrow \infty} \frac{\kappa^L(x, x')}{c(L)} = c_1 > 0 \text{ for any } x \neq x', \text{ and } \lim_{L \rightarrow \infty} \frac{\kappa^L(x, x)}{c(L)} \neq c_1, \quad (27)$$

then

$$\lim_{L \rightarrow \infty} f(x) = \operatorname{sgn} \sum_{i=1}^n y^i \quad \text{MAJORITY CLASSIFIER} \quad (28)$$

► What ϕ 's satisfy theorem? ReLU, all ϕ with $B \neq 1$.



Case $A = A' = 0$: Networks implement 1-nearest neighbor

Theorem (Theorem 3 in [Radhakrishnan, Belkin, Ulher, 2022])

Given x , assume w.l.o.g. that $x^T x^1 = \max_{i=1:n} x^T x^i$.

$$\lim_{L \rightarrow \infty} \frac{\kappa^L(x, x^i)}{\kappa^L(x, x^1)} = 0. \quad (29)$$

and

$$\lim_{L \rightarrow \infty} f(x) = \text{sgn} y^1 \quad 1\text{-nn} \quad (30)$$



Case $A = 0$, $A' \neq 0$: Networks implement singular kernel classifier

Theorem (Theorem 1 in [Radhakrishnan, Belkin, Usher, 2022])

$$\lim_{L \rightarrow \infty} \frac{\kappa^L(x, x')}{(A')^{2L}(L+1)} = \frac{R(\|x - x'\|)}{\|x - x'\|^\alpha}, \quad (31)$$

with $\alpha = -4 \frac{\log A'}{\log B'}$ and $R() \geq 0$, bounded, and $R(u) > \delta$ around 0.

- ▶ if $\alpha > 0$, $\frac{R(\|x - x'\|)}{\|x - x'\|^\alpha}$ is **singular kernel**
- ▶ Computationally not a problem: if data $x^{1:n}$ distinct, G_0 is well defined
- ▶ If $x = x^i$, set $f(x) = y^i$.

Optimality of singular kernel classifier

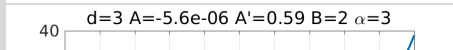
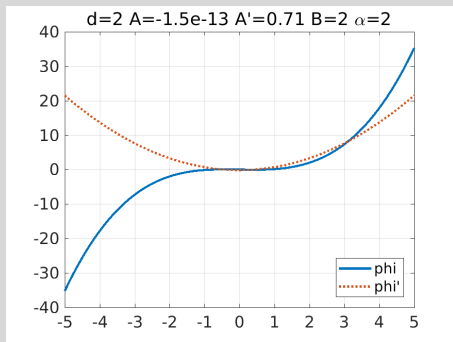
Theorem (Theorem 2 in [Radhakrishnan, Belkin, Ulher, 2022])

If $A = 0$, $A' \neq 0$ and $\alpha = \frac{1}{d}$ then $\lim_{L \rightarrow \infty} f(x)$ is Bayes-optimal.

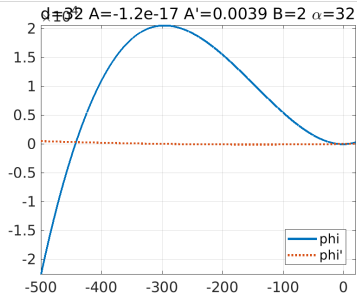
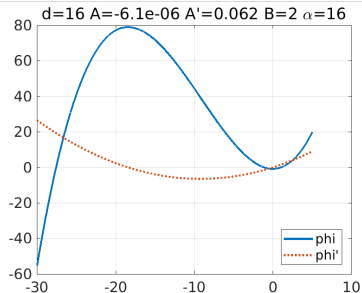
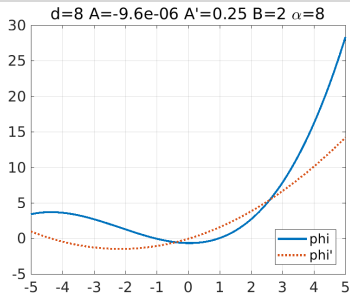
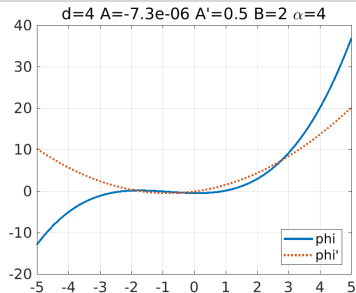
► What activations ϕ satisfy this theorem?

$$\phi^{\text{opt}}(z) = \frac{1}{2^{d/4}} \frac{z^3 - 3z}{\sqrt{6}} + \sqrt{1 - 2^{1-d/2}} \frac{z^2 - 1}{\sqrt{2}} + \frac{1}{2^{d/4}} z \quad \text{for } d \geq 2$$

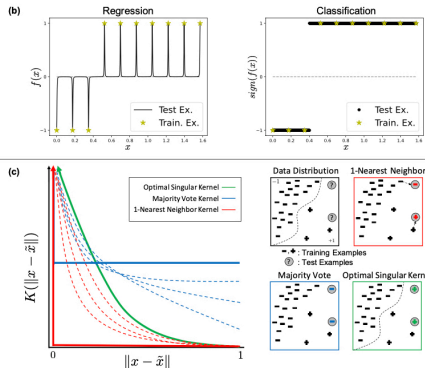
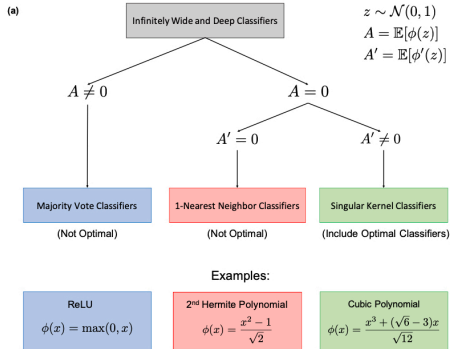
. (32)



Optimal singular kernels for $d = 4, 8, 16, 32$



Summary



- b** when $A \neq 0$, $\lim_{L \rightarrow \infty} \kappa^L(x, x') = 0$ for $x \neq x'$, and $f(x) = 0$ is vanishingly small (useless for regression), but $\text{sgn}f(x)$ can be optimal for classification
- c** Singular kernel $\alpha > d$, $\alpha < d$, majority vote kernel, and 1-nn kernel

Limits of some activation functions

ϕ^{opt} Bayes classifier

ReLU majority vote

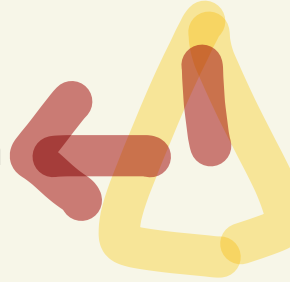
sigmoid $\frac{1}{1+e^{-z}} - \frac{1}{2}$ 1-nearest neighbor

Back of Form; Q1

TAs ability to help you when you interacted with him

A=poor.....E=excellen

Please share any feedback for the TA



Back of yellow form