

Neural networks, wide and deep, singular kernels and Bayes optimality

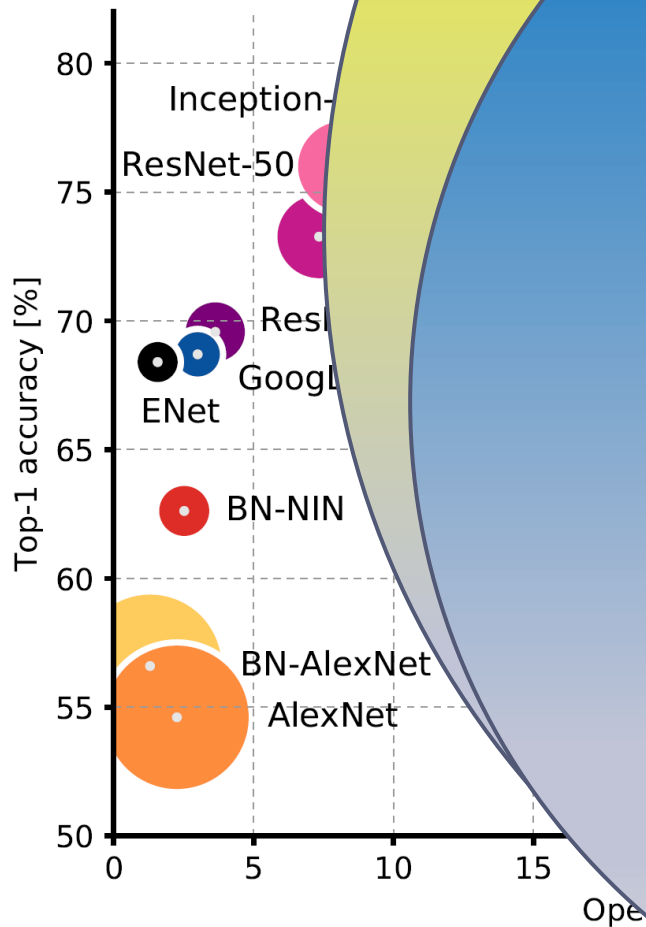
Mikhail Belkin

University of California San Diego,
Halıcıoğlu Data Science Institute

Collaborators: Adit Radhakrishnan, Caroline Uhler, Chaoyue Liu, Libin Zhu, Like Hui, Alexander Rakhlin, Alexander Tsybakov

Non-Linear and High Dimensional Inference,
IHP, Paris, Oct 2022

L1



Switch Transformer, 2021:
1.6 trillion parameters

From Canziani, et al., 2017.

This talk

- ▶ A taxonomy and of infinitely wide and deep interpolating networks.
- ▶ Equivalence with singular kernel machines for classification.
- ▶ Bayes optimality for classification but **not regression**.



The problem of machine Learning

Input: data (x_i, y_i) , $i = 1..n$, $x_i \in \mathbb{R}^d$, $y_i \in \{-1, 1\}$ (classification)

Goal: construct $f^*: \mathbb{R}^d \rightarrow \mathbb{R}$, that best “generalizes” to new data.


Under the standard statistical assumptions:

$$f^* = \underset{f}{\operatorname{argmin}} E_{\text{unseen data}} L(f(x), y)$$



Empirical Risk Minimization

Most ML algorithms are based on minimizing empirical risk.



Empirical risk

$$f_{ERM}^* = \operatorname{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum_{\text{training data}} L(f(x_i), y_i)$$

Tension between empirical risk and expected risk.



Interpolation and benign over-fitting

Recent understanding: interpolation does not contradict generalization.

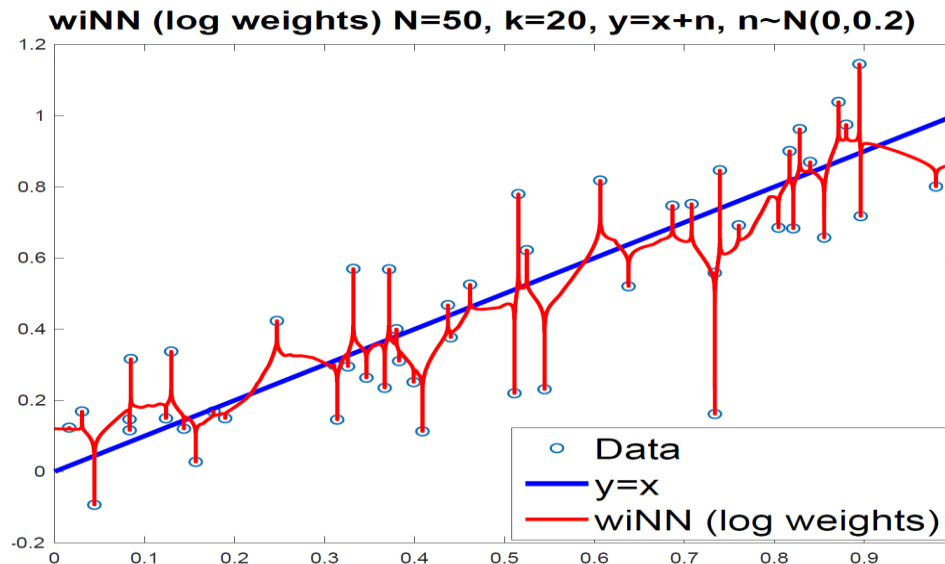


Table 18: ASR results on training and test set, error rate

Model	Task	train with square loss (%)		train with cross-entropy (%)		square loss w/ same epochs as CE (%)	
		Train	Test	Train	Test	Train	Test
Attention+CTC	TIMIT (PER)	0.9	20.8	4.8	20.8	0.9	20.8
(Kim et al., 2017)	TIMIT (CER)	4.5	32.5	11.6	33.4	4.5	32.5
VGG+BLSTMP	WSJ (WER)*	0.7	5.1	0.3	5.3	0.7	5.1
(Moritz et al., 2019)	WSJ (CER)*	0.3	2.4	0.1	2.5	0.3	2.4
VGG+BLSTM	Librispeech (WER)*	0.8	9.8	0.4	10.6	0.8	10.3
(Moritz et al., 2019)	Librispeech (CER)*	0.6	9.7	0.3	10.7	0.6	10.2
Transformer	WSJ (WER)*	0.7	5.7	0.5	5.8	0.7	5.7
(Watanabe et al., 2018)	Librispeech (WER)*	0.9	9.4	1.2	9.2	0.9	9.4

* For WSJ and Librispeech, we take 10% of the training set for the evaluation of the training error rate.

Train loss at **optimal early stopping** is far below test loss.

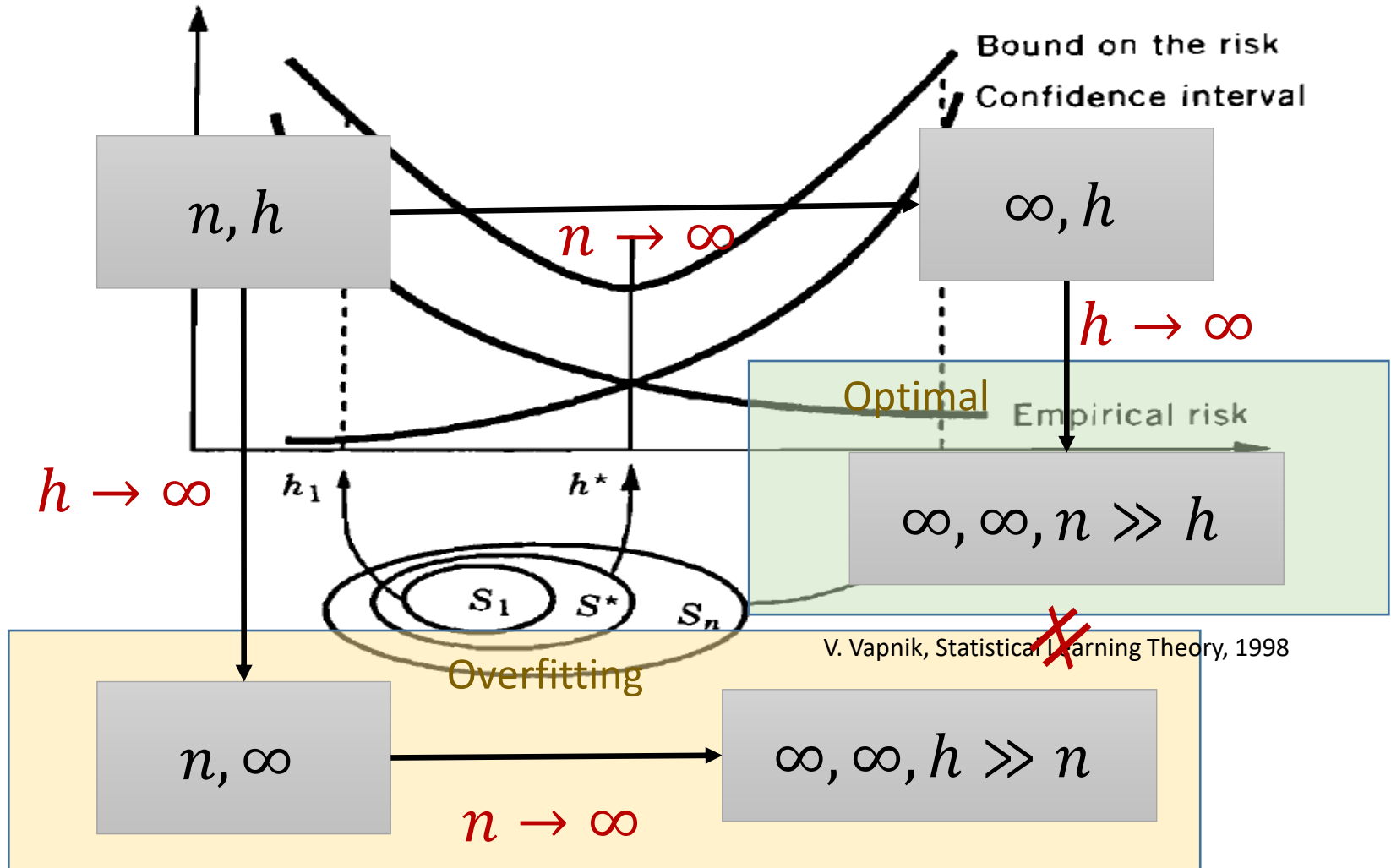
[Hui, B., ICLR 2021]



Classical: non-commutative

number of points n
hypothesis class h

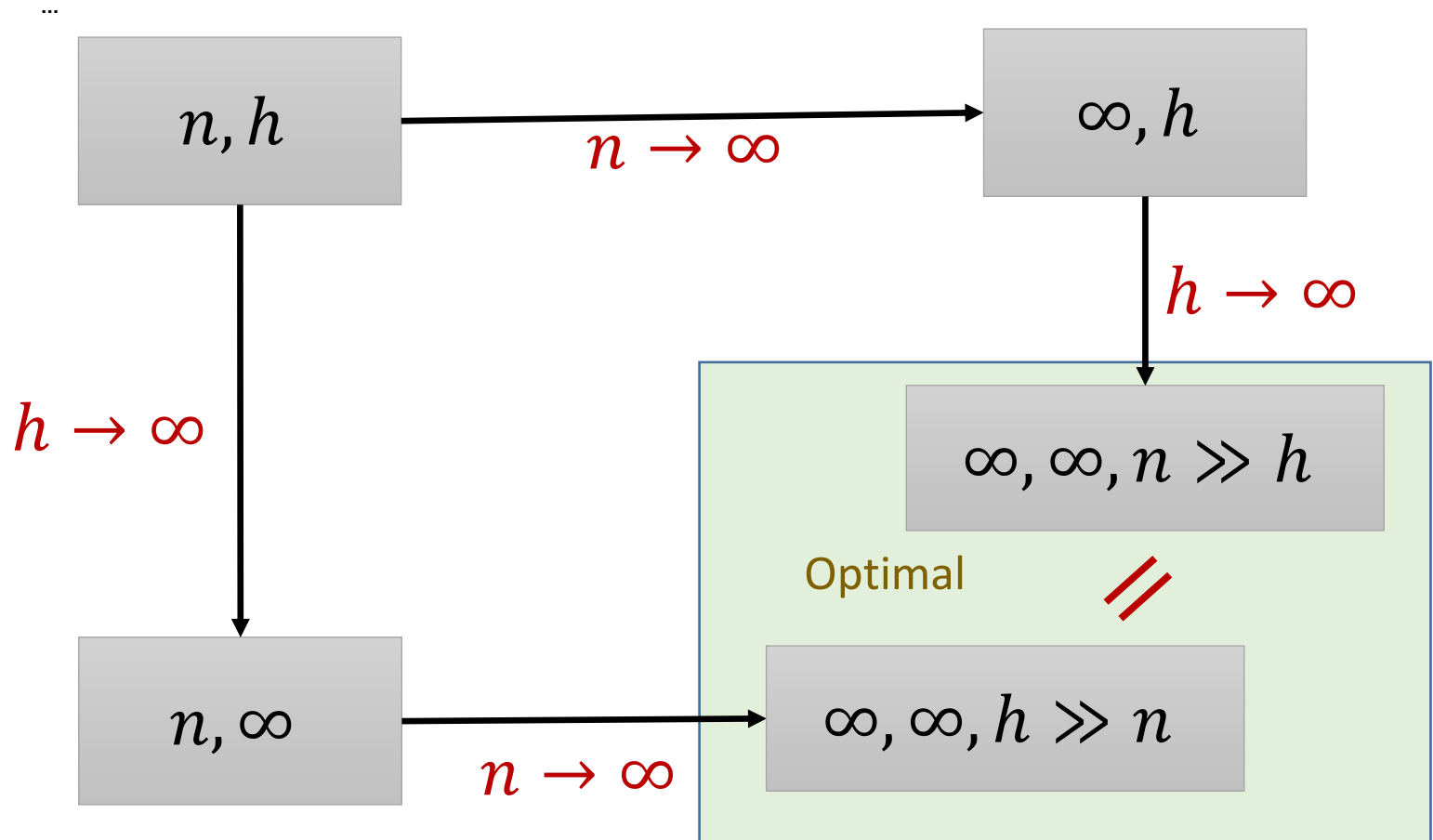
6.1 THE SCHEME OF THE STRUCTURAL RISK MINIMIZATION INDUCTION PRINCIPLE



“Modern” commutative

number of points n

hypothesis class h



kernel machines

Beautiful classical statistical/mathematical theory based on Reproducing Kernel Hilbert Spaces (RKHS) -- Hilbert Space of functions with bounded evaluation functionals.

RKHS Theory [Aronszajn, ..., 1950s]

Splines [Parzen, Wahba, ..., 1970-80s]

Kernel machines [Vapnik, ..., 1990s]

Wide neural networks [Jacot, Gabriel, Hogler, ..., 2020s]



kernels

Any RKHS corresponds to a PSD kernel.

$$k(x, y) = e^{-\frac{\|x-y\|^2}{\sigma^2}}$$

Gaussian kernel

$$k(x, y) = e^{-\frac{\|x-y\|}{\sigma}}$$

Laplace kernel

Many others...



Interpolating kernel machines

$$f^* = \operatorname{argmin}_{f \in \mathcal{H}, \forall_i f(x_i) = y_i} \|f\|_{\mathcal{H}}$$

Representer theorem -- solution:

$$f^*(x) = \sum_i \alpha_i k(x_i, x), \quad \boldsymbol{\alpha} = K^{-1} \mathbf{y}$$

$$K_{ij} = k(x_i, x_j)$$



width -- Transition to linearity

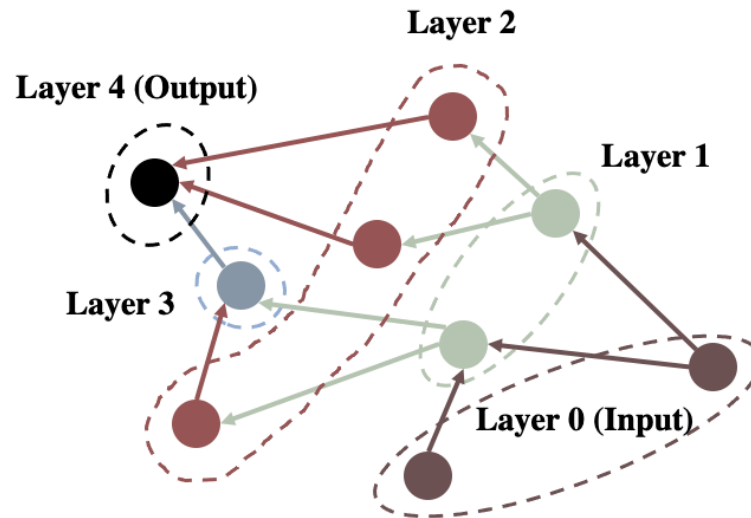
very wide neural networks (w. linear output layer)
= linear functions of parameters
= kernel machines.

$$k_w(x, z) = \langle \nabla_w f_w(x), \nabla_w f_w(z) \rangle$$

First identified in [Jacot, Gabriel, Hogler, 18] as constant NTK along the training trajectory for the model $f_w(x)$.



Transition to linearity in DAG neural networks



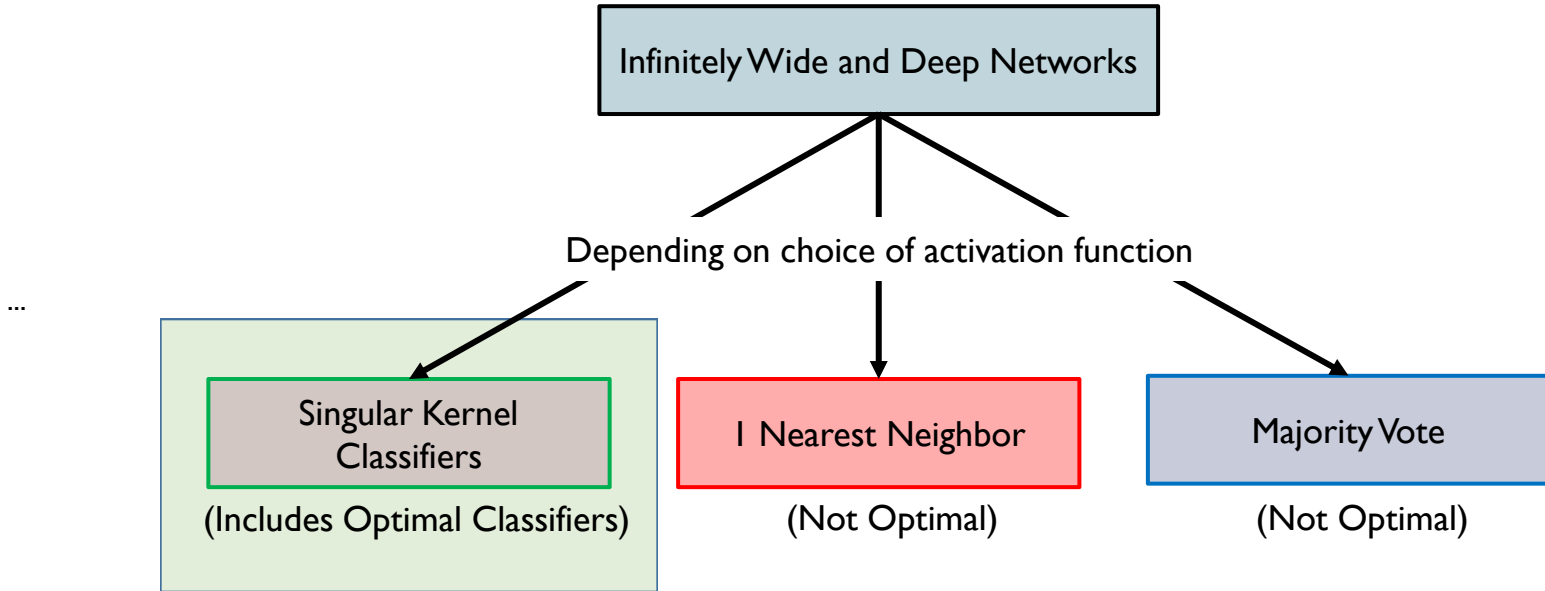
Theorem:

For a neural network (w. arbitrary activations) corresponding to a DAG with minimum in-degree m (width) and linear output layer

$$\|H(F)(w)\| = O(1/\sqrt{m})$$

[Zhu, Liu, B., NeurIPS 2022]

Width + Depth



[Radhakrishnan, B., Uhler, 2022]

Inverse and Direct methods

Kernel machine (inverse):

$$y(x) = \text{sign} (Y_n K_n^{-1} K(X, x))$$

$$(K_n)_{ij} = k(x_i, x_j) \quad Y_n = (y_1, \dots, y_n),$$

$$K(X, x) = (k(x_1, x), \dots, k(x_n, x))^T$$

Kernel smoother (Nadaraya-Watson) (direct):

$$y(x) = \text{sign} \left(\frac{\sum y_i k(x_i, x)}{\sum k(x_i, x)} \right) = \text{sign} (Y_n K(X, x))$$

Singular kernels and Interpolating Nadaraya-Watson schemes

$$f(x) = \frac{\sum y_i k(x_i, x)}{\sum k(x_i, x)}$$

$$k(x_i, x) = \frac{1}{\|x - x_i\|^\alpha}$$

Claim: NW predictors with singular kernels are interpolating schemes: $\forall_i f(x) = y_i$

(Shepard's interpolation, 1968)

Optimality of Interpolating Nadaraya-Watson schemes

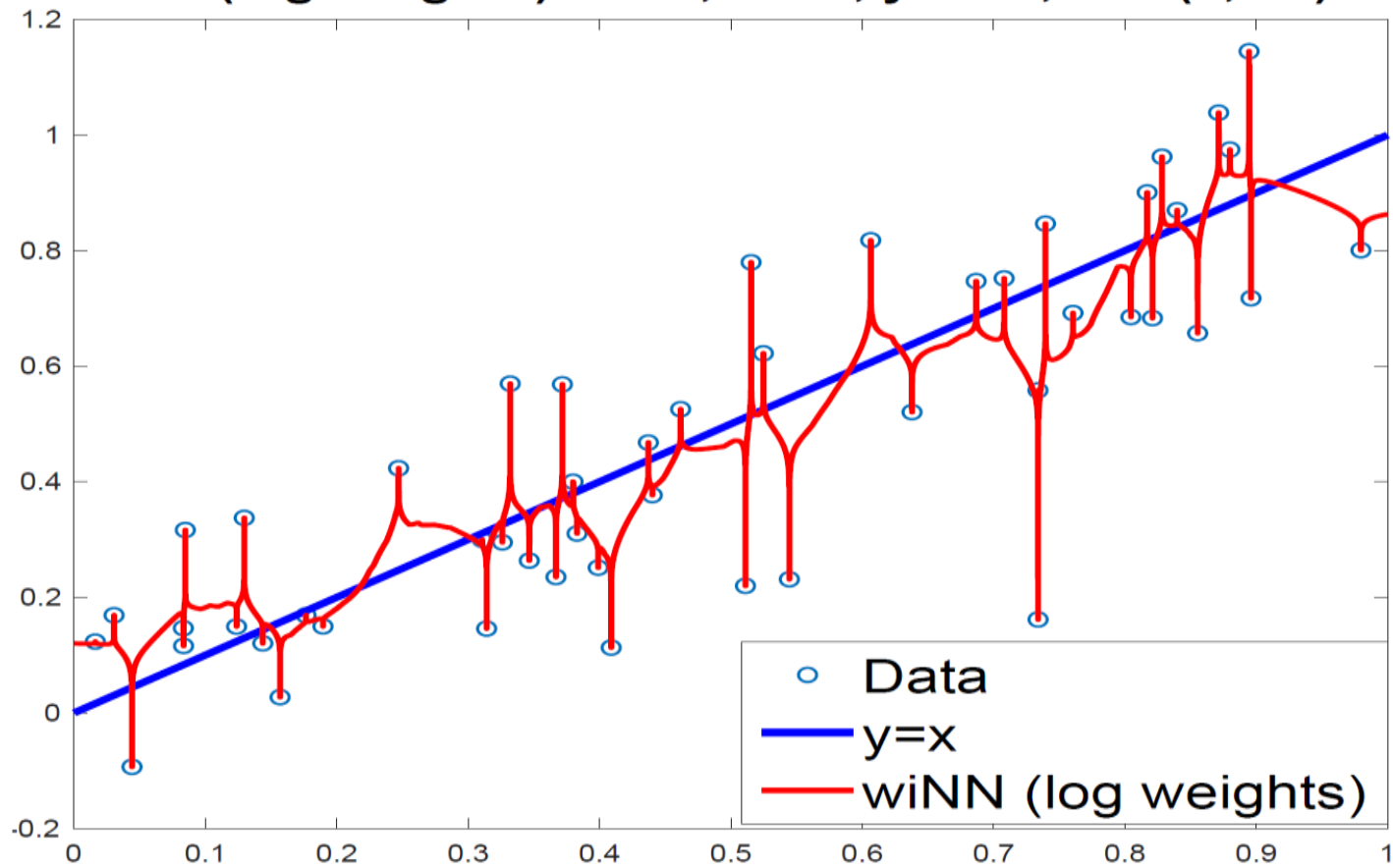
Consistency for regression. [Devroye, Györfi, Krzyzak, 98]
(the Hilbert scheme) $\alpha = d$.

Optimality for regression/classification:

Weighted interpolated schemes with certain singular kernels are consistent (converge to Bayes optimal) for classification in any dimension. Moreover, **statistically (minimax) optimal** for regression in any dimension.

[B., Hsu, Mitra, NeurIPS 18], followup [B., Rakhlin, Tsybakov, AISTATS 19]

wiNN (log weights) N=50, k=20, $y=x+n$, $n\sim N(0,0.2)$



Equivalence of direct and inverse methods

Claim. Direct and inverse methods with singular kernels are equivalent for classification:

$$y(x) = \text{sign} \left(\frac{\sum y_i k(x_i, x)}{\sum k(x_i, x)} \right) = \text{sign} (Y_n K_n^{-1} K(X, x))$$

Proof.

1. $\text{sign} \left(\frac{\sum y_i k(x_i, x)}{\sum k(x_i, x)} \right) = \text{sign} (Y_n K(X, x))$

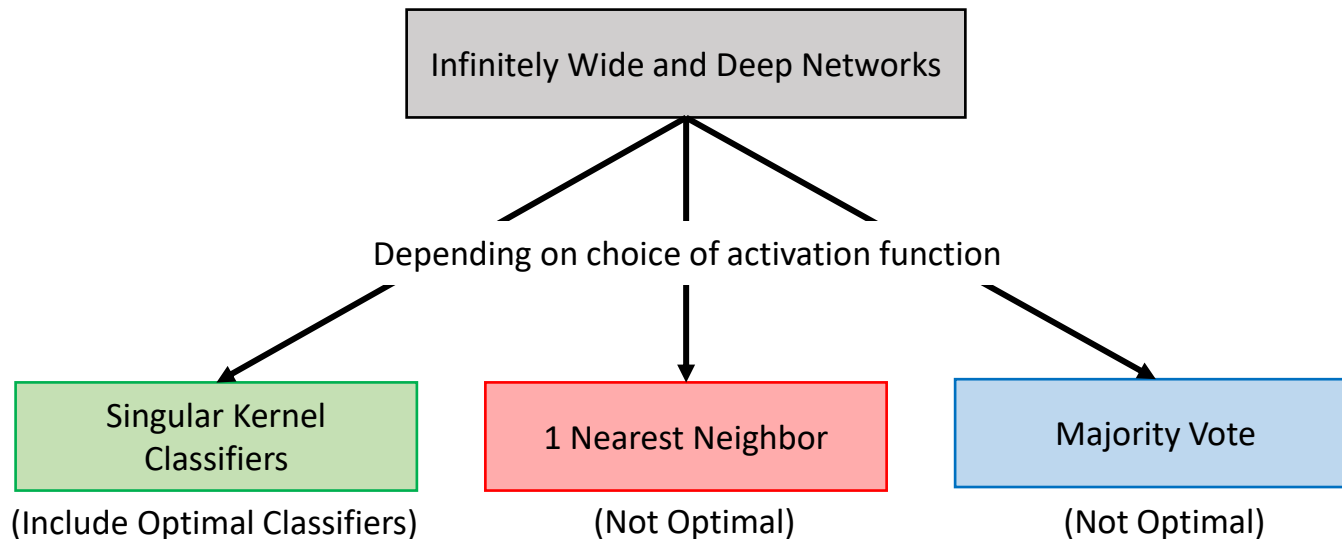
Off-diagonal terms

2. singular kernels: $K_n = \infty I + G_n$

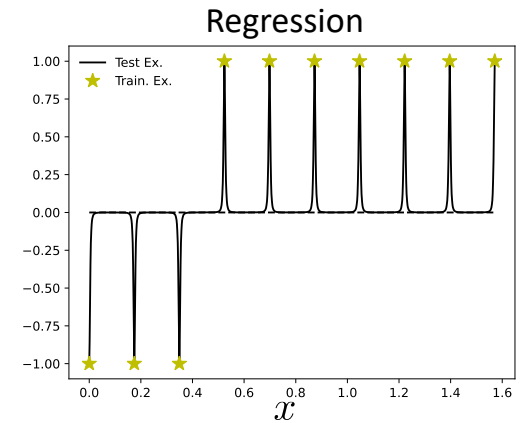
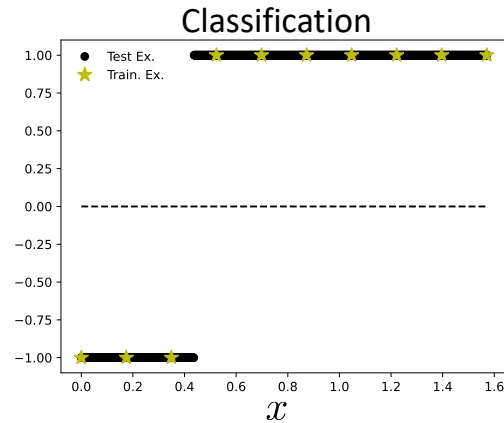
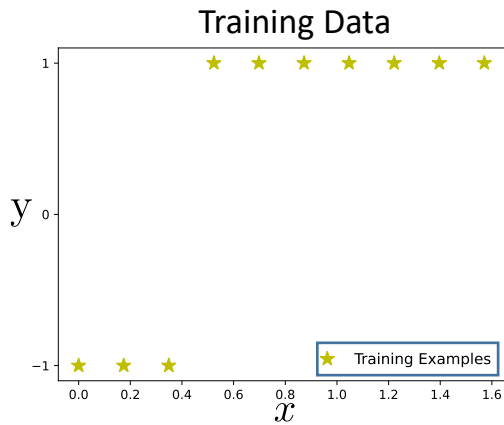
$$K_n^{-1} = (\infty I + G_n)^{-1} = \frac{1}{\infty} \left(I + \frac{G_n}{\infty} \right)^{-1} = \frac{1}{\infty} \left(I - \frac{1}{\infty} G_n \right) = \frac{1}{\infty} I$$

wide and deep networks

- A taxonomy of infinitely wide and deep classifiers (on a sphere) based on activation function.
- Construct infinitely wide and deep FC networks that, when trained using standard methods, achieve optimality for classification.



Regression vs Classification



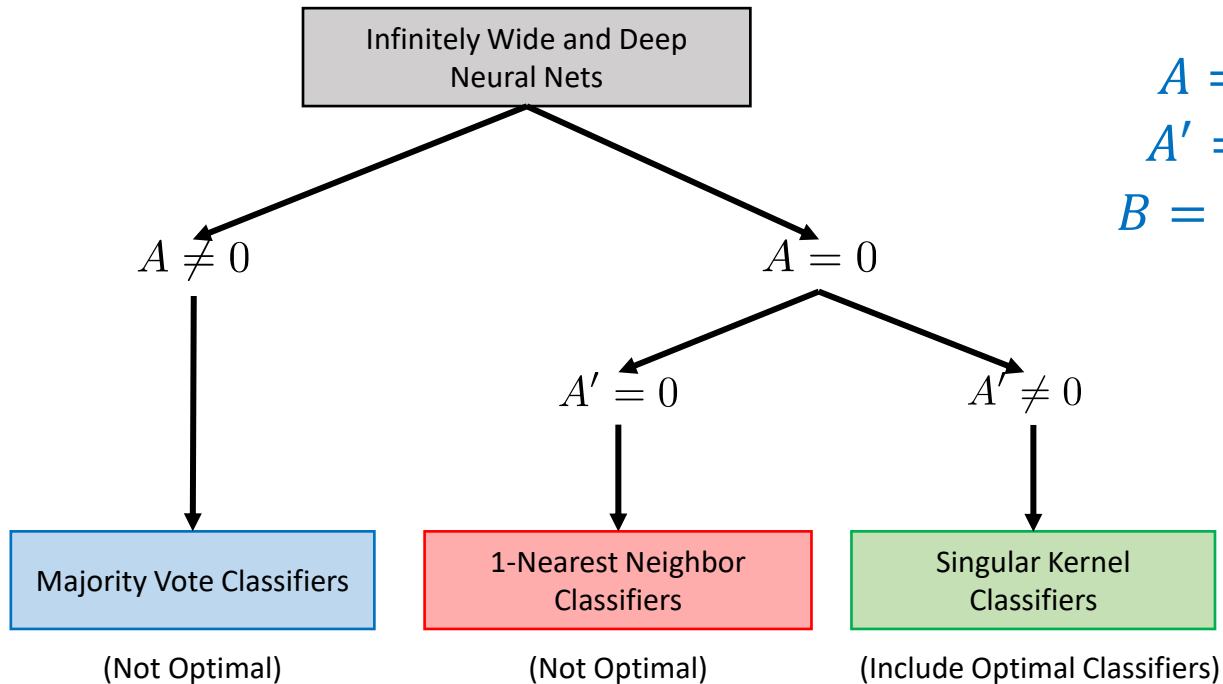
- Converges to delta function for regression but not for classification!

Taxonomy of Infinitely wide and Deep FC Nets

Depends on activation function ϕ .

$$z \sim \mathcal{N}(0,1)$$

$$A = \mathbb{E}(\phi(z))$$
$$A' = \mathbb{E}(\phi'(z))$$
$$B = \mathbb{E}([\phi'(z)]^2)$$



Examples:

ReLU

$$\phi(x) = \max(0, x)$$

2nd Hermite Polynomial

$$\phi(x) = \frac{x^2 - 1}{\sqrt{2}}$$

Cubic Polynomial

$$\phi(x) = \frac{x^3 + (\sqrt{6} - 3)x}{\sqrt{12}}$$

Singular Kernel Classifiers

Theorem A. ($A = 0, A' \neq 0$)

$$k^\infty(x, z) \sim \left(\frac{R(\|x-z\|)}{\|x-z\|^\alpha} \right)$$

Infinite depth NTK

$$\alpha = -\frac{4 \ln A'}{\ln B}$$

Consequence (using Devroye, Györfi, and Krzyżak (1998)):

If $\alpha = d$, infinite depth neural net predictor h^∞ is Bayes optimal.

1-NN and Majority vote

Theorem B. ($A = 0 = A' = 0$)

$$\text{sign}(h^\infty) = 1 - NN(x)$$

Theorem C*. ($A \neq 0$)

$$\text{sign}(h^\infty) = \text{majority vote } (y_1, \dots, y_n)$$

what is going on?

Kernels for deep FC networks (on the unit sphere). Put $v = \langle x, z \rangle$

$$k^0(v) = v$$

$$k^l(v) = \underbrace{\psi(\dots(\psi(v)\dots))}_l + k^{l-1}(v) \underbrace{\psi'(\psi(\dots(\psi(v)\dots))}_{l-1}$$

$$\psi(v) = \mathbb{E}_{(u,t) \sim \mathcal{N}(0, \Lambda(v))} \phi(u)\phi(t), \quad \Lambda = \begin{pmatrix} 1 & v \\ v & 1 \end{pmatrix}$$

Dual activation function.

First key observation

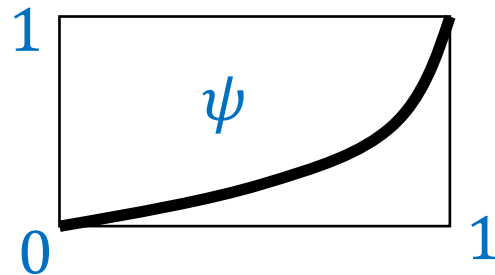
$$k^l(v) = \underbrace{\psi(\dots(\psi(v)\dots))}_l + k^{l-1}(v)\psi'(\psi(\dots(\psi(v)\dots)))$$

$$\psi^l(v) = \psi(\dots(\psi(v)\dots))$$

kernels $k^l(\langle x, z \rangle)$ and $\psi^l(\langle x, z \rangle)$ (after normalization) have poles of the same order at $v = 1$ as $l \rightarrow \infty$.

Iteration

$\psi: [0,1] \rightarrow [0,1]$ is **convex** and **monotone**.



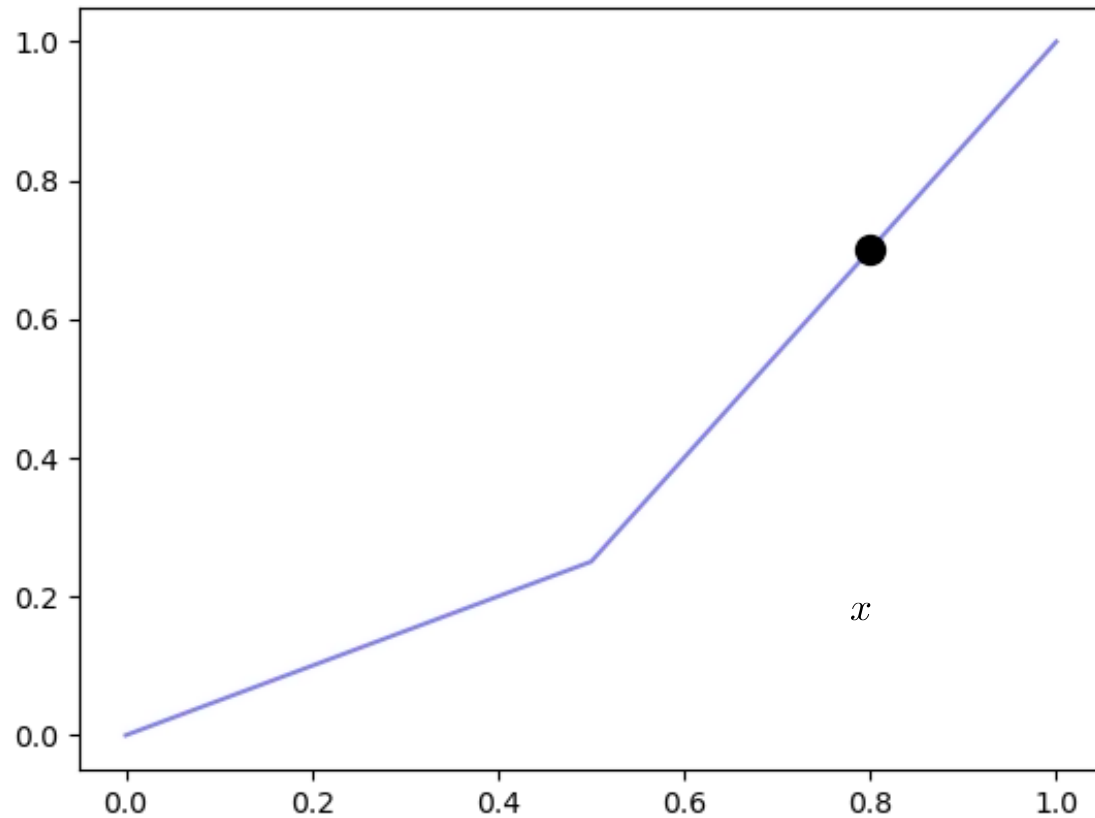
Two fixed points (case 1, sing. kernel):

$$A = \psi(0) = 0, \quad A' = \psi(1) = 1$$

What does $\psi^l(v) = \psi(\dots(\psi(v)\dots))$ look like?

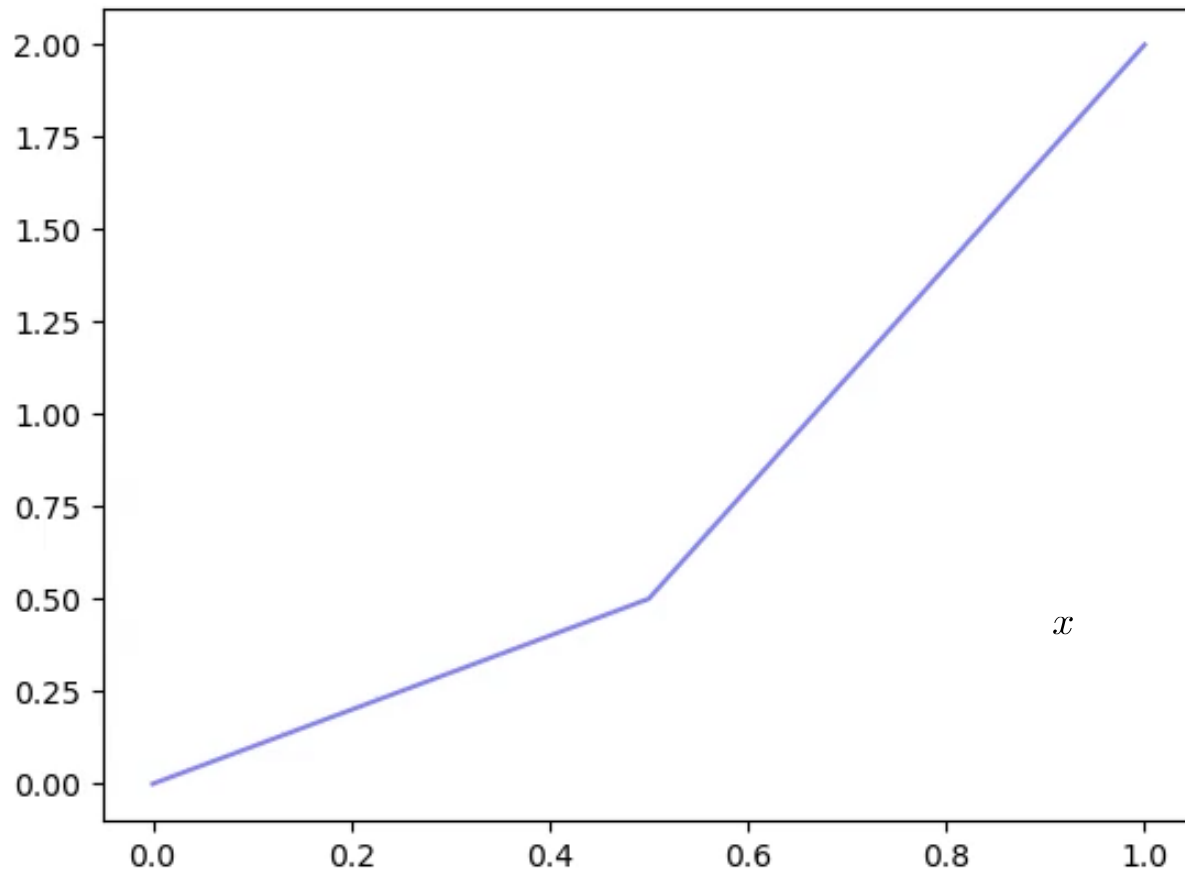
Piecewise linear function

Iteration of Function (Single Point)

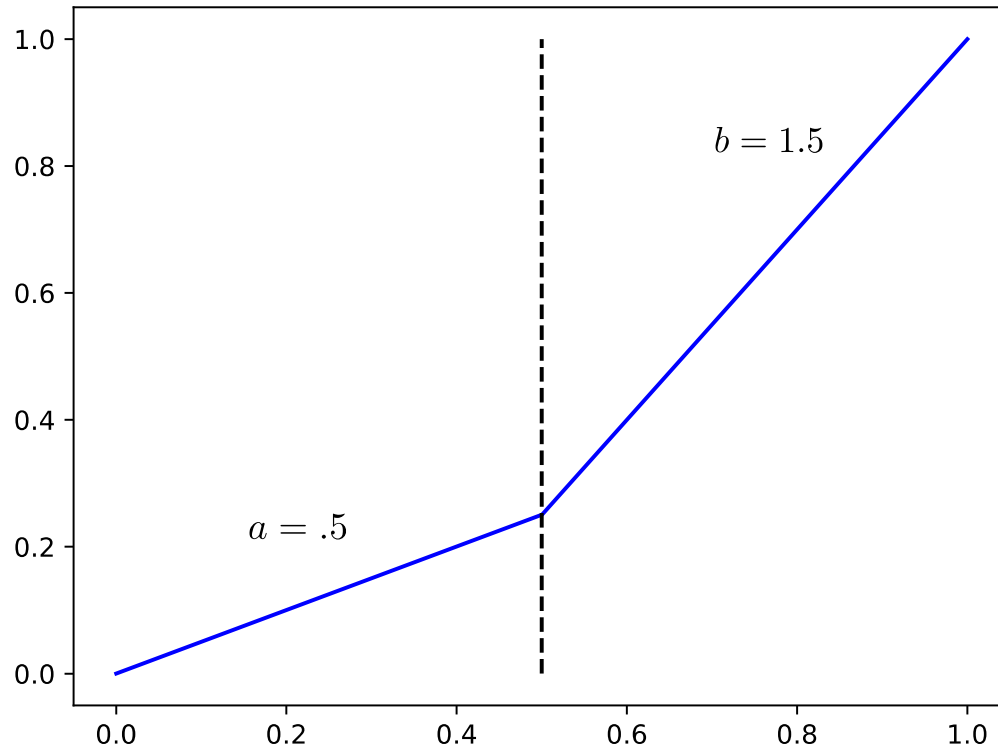


Iteration on the unit interval

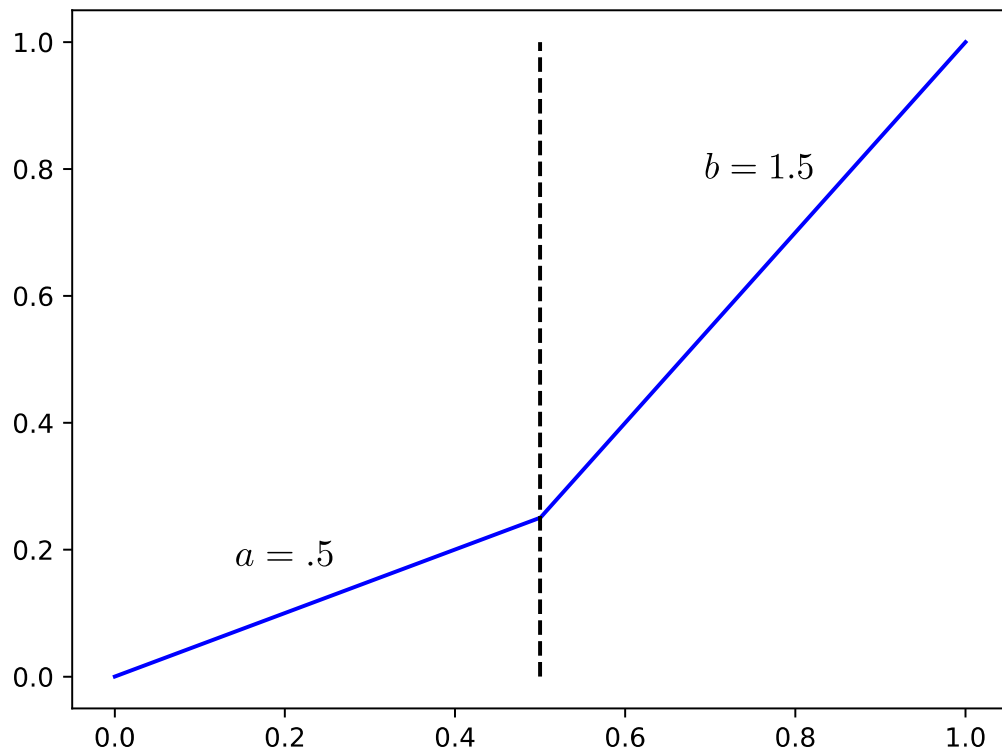
Iteration of Normalized Function



Iteration



$\psi^\infty(v) = \lim_{l \rightarrow \infty} \frac{\psi^l(v)}{a^l}$ is well-defined and has a pole at 1.

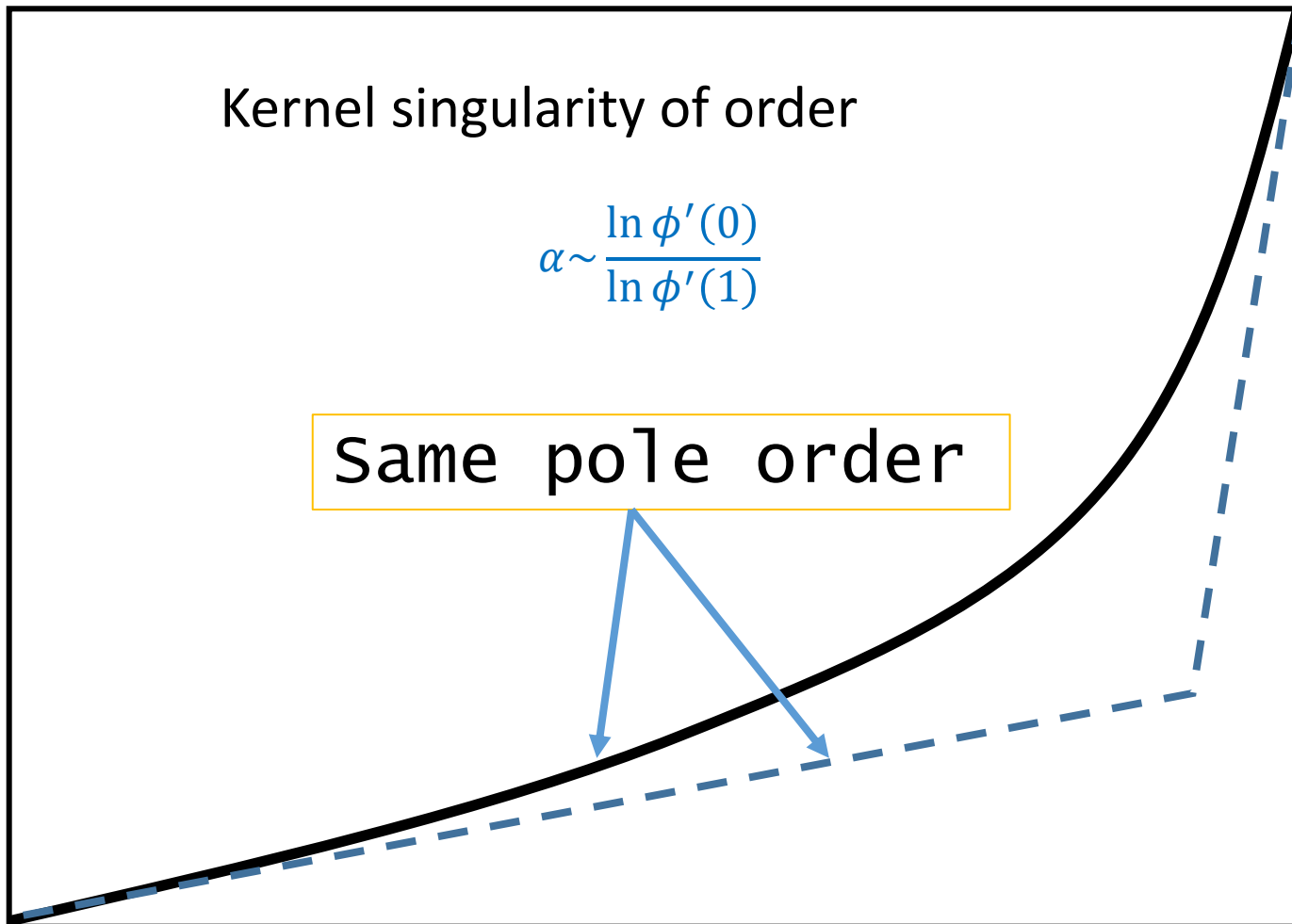


What matters are the slopes.

Kernel singularity of order

$$\alpha \sim -\frac{\ln b}{\ln a}$$

Key observation



Summary

Remarkable properties of neural networks:

- wide neural networks = kernel machines.
- Deep neural networks = kernel machines with singular kernels = NW schemes for classification.
- Clear separation between regression and classification.