AND THE

CURSE OF DIMENSIONALITY

arXiv:1905.10843

Stefano Spigler

Jonas Paccolat, Mario Geiger, Matthieu Wyart



SUPERVISED DEEP LEARNING

• Why and how does deep **supervised** learning work?

• Learn from examples: **how many** are needed?

- Typical tasks:
 - Regression (fitting functions)
 - Classification

LEARNING CURVES

• Performance is evaluated through the **generalization error** ϵ

• Learning curves decay with number of examples *n*, often as

 $\epsilon \sim n^{-eta}$

• β depends on the **dataset** and on the **algorithm**

Deep networks: $eta \sim 0.07$ -0.35 [Hestness et al. 2017]

We lack a theory for β for deep networks!

LINK WITH KERNEL LEARNING

• Performance increases with **overparametrization**

[Neyshabur et al. 2017, 2018, Advani and Saxe 2017] [Spigler et al. 2018, Geiger et al. 2019, Belkin et al. 2019]

 \longrightarrow study the infinite-width limit!

[Mei et al. 2017, Rotskoff and Vanden-Eijnden 2018, Jacot et al. 2018, Chizat and Bach 2018, ...]



LINK WITH KERNEL LEARNING

• Performance increases with **overparametrization**

[Neyshabur et al. 2017, 2018, Advani and Saxe 2017] [Spigler et al. 2018, Geiger et al. 2019, Belkin et al. 2019]

 \longrightarrow study the infinite-width limit!

[Mei et al. 2017, Rotskoff and Vanden-Eijnden 2018, Jacot et al. 2018, Chizat and Bach 2018, ...]



• With a specific scaling, infinite-width limit \rightarrow **kernel learning** [Jacot et al. 2018] (next slides) Neural Tangent Kernel

What are the learning curves of kernels like?

OUTLINE

• Very brief introduction to kernel methods

• Performance of kernels on real data

• Gaussian data: Teacher-Student regression

• Gaussian approximation: smoothness and effective dimension

• Dimensional reduction via invariants in the task

- Kernel methods learn non-linear functions or boundaries
- Map data to a **feature space**, where the problem is linear

data $\underline{x} \longrightarrow \phi(\underline{x}) \longrightarrow$ use linear combination of features



- Kernel methods learn non-linear functions or boundaries
- Map data to a **feature space**, where the problem is linear

data $\underline{x} \longrightarrow \phi(\underline{x}) \longrightarrow$ use linear combination of features





kernel $K(\underline{x}, \underline{x}')$

- Kernel methods learn non-linear functions or boundaries
- Map data to a **feature space**, where the problem is linear

data $\underline{x} \longrightarrow \phi(\underline{x}) \longrightarrow$ use linear combination of features



E.g. kernel regression:

• Target function $\underline{x}_{\mu} o Z(\underline{x}_{\mu}), \ \mu = 1, \dots, n$

E.g. kernel regression:

• Target function $\underline{x}_{\mu} o Z(\underline{x}_{\mu}), \ \mu = 1, \dots, n$

• Build an estimator $\hat{Z}_K(\underline{x}) = \sum_{\mu=1}^n c_\mu K(\underline{x}_\mu, \underline{x})$

E.g. kernel regression:

• Target function $\underline{x}_{\mu} o Z(\underline{x}_{\mu}), \ \mu = 1, \dots, n$

• Build an estimator $\hat{Z}_K(\underline{x}) = \sum_{\mu=1}^n c_\mu K(\underline{x}_\mu, \underline{x})$

• Minimize training MSE
$$= rac{1}{n} \sum_{\mu=1}^n \left[\hat{Z}_K(\underline{x}_\mu) - Z(\underline{x}_\mu)
ight]^2$$

E.g. kernel regression:

• Target function $\underline{x}_{\mu} o Z(\underline{x}_{\mu}), \ \mu = 1, \dots, n$

• Build an estimator $\hat{Z}_K(\underline{x}) = \sum_{\mu=1}^n c_\mu K(\underline{x}_\mu, \underline{x})$

• Minimize training MSE = $\frac{1}{n} \sum_{\mu=1}^{n} \left[\hat{Z}_{K}(\underline{x}_{\mu}) - Z(\underline{x}_{\mu}) \right]^{2}$

• Estimate the **generalization error** $\epsilon = \mathbb{E}_{\underline{x}} \left[\hat{Z}_K(\underline{x}) - Z(\underline{x}) \right]^2$

REPRODUCING KERNEL HILBERT SPACE (RKHS)

A kernel K induces a corresponding Hilbert space \mathcal{H}_K with norm

$$\|Z\|_K = \int \mathrm{d}\underline{x}\mathrm{d}\underline{y}\,Z(\underline{x})K^{-1}(\underline{x},\underline{y})Z(\underline{y})$$

where $K^{-1}(\underline{x}, y)$ is such that

$$\int \mathrm{d} \underline{y} \, K^{-1}(\underline{x}, \underline{y}) K(\underline{y}, \underline{z}) = \delta(\underline{x}, \underline{z})$$

 \mathcal{H}_K is called the **Reproducing Kernel Hilbert Space** (RKHS)

PREVIOUS WORKS

Regression: performance depends on the target function!

PREVIOUS WORKS

Regression: performance depends on the target function!

• If only assumed to be Lipschitz, then
$$\beta = \frac{1}{d}$$

Curse of dimensionality! [Luxburg and Bousquet 2004]

PREVIOUS WORKS

Regression: **performance depends on the target function!**

• If only assumed to be **Lipschitz**, then
$$\beta = \frac{1}{d}$$

Curse of dimensionality! [Luxburg and Bousquet 2004]

• If assumed to be in the **RKHS**, then $\beta \geq \frac{1}{2}$ does not depend on d

[Smola et al. 1998, Rudi and Rosasco 2017]

Regression: performance depends on the target function!

• If only assumed to be **Lipschitz**, then
$$\beta = \frac{1}{d}$$

Curse of dimensionality! [Luxburg and Bousquet 2004]

• If assumed to be in the **RKHS**, then $\beta \geq \frac{1}{2}$ does not depend on d

[Smola et al. 1998, Rudi and Rosasco 2017]

• Yet, RKHS is a very strong assumption on the smoothness of the target function (see later on)

[Bach 2017]

REAL DATA AND ALGORITHMS

We apply kernel methods on





- Same exponent for regression and classification
- Same exponent for Gaussian and Laplace kernel
- MNIST and CIFAR10 display exponents $\beta \gg \frac{1}{d}$ but $< \frac{1}{2}$



- Same exponent for regression and classification
- Same exponent for Gaussian and Laplace kernel
- MNIST and CIFAR10 display exponents $\beta \gg \frac{1}{d}$ but $< \frac{1}{2}$

• Controlled setting: Teacher-Student regression

• Controlled setting: **Teacher-Student regression**

• Training data are sampled from a Gaussian Process:

 $Z_T(\underline{x}_1), \ldots, Z_T(\underline{x}_n) \sim \mathcal{N}(0, K_T)$ \underline{x}_{μ} are random on a *d*-dim hypersphere

• Controlled setting: **Teacher-Student regression**

• Training data are sampled from a Gaussian Process:

 $Z_T(\underline{x}_1), \dots, Z_T(\underline{x}_n) \sim \mathcal{N}(0, K_T)$ $\underline{x}_{\mu} \text{ are random on a } \boldsymbol{d}\text{-dim hypersphere}$ $\mathbb{E}Z_T(\underline{x}_{\mu}) = 0$

 $\mathbb{E} Z_T(\underline{x}_\mu) Z_T(\underline{x}_
u) = K_T(\|\underline{x}_\mu - \underline{x}_
u\|)$

• Controlled setting: **Teacher-Student regression**

• Training data are sampled from a Gaussian Process:

 $Z_T(\underline{x}_1), \dots, Z_T(\underline{x}_n) \sim \mathcal{N}(0, K_T)$ $\underline{x}_{\mu} \text{ are random on a } \boldsymbol{d}\text{-dim hypersphere}$ $\mathbb{E}Z_T(\underline{x}_{\mu}) Z_T(\underline{x}_{\nu}) = K_T(||\underline{x}_{\mu} - \underline{x}_{\nu}||)$

• Regression is done with another kernel K_S

TEACHER-STUDENT: SIMULATIONS



$$\frac{Regression}{Regression}: \frac{\hat{Z}_{S}(\underline{x}) = \sum_{\mu=1}^{n} c_{\mu} K_{S}(\underline{x}_{\mu}, \underline{x})}{\text{Minimize} = \frac{1}{n} \sum_{\mu=1}^{n} \left[\hat{Z}_{S}(\underline{x}_{\mu}) - Z_{T}(\underline{x}_{\mu}) \right]^{2}}$$

$$\hat{Z}_{S}(\underline{x}) = \underline{k}_{S}(\underline{x}) \cdot \mathbb{K}_{S}^{-1} \underline{Z}$$

where

$$\frac{\hat{Z}_{S}(\underline{x}) = \sum_{\mu=1}^{n} c_{\mu} K_{S}(\underline{x}_{\mu}, \underline{x})}{\text{Minimize} = \frac{1}{n} \sum_{\mu=1}^{n} \left[\hat{Z}_{S}(\underline{x}_{\mu}) - Z_{T}(\underline{x}_{\mu}) \right]^{2}}$$

Explicit solution:

$$\hat{Z}_{S}(\underline{x}) = \underline{k}_{S}(\underline{x}) \cdot \mathbb{K}_{S}^{-1} \underline{Z}_{T}$$

where

$$(\underline{k}_S(\underline{x}))_\mu = K_S(\underline{x}_\mu, \underline{x})$$
kernel overlap

$$\frac{\hat{Z}_{S}(\underline{x}) = \sum_{\mu=1}^{n} c_{\mu} K_{S}(\underline{x}_{\mu}, \underline{x})}{\text{Minimize} = \frac{1}{n} \sum_{\mu=1}^{n} \left[\hat{Z}_{S}(\underline{x}_{\mu}) - Z_{T}(\underline{x}_{\mu}) \right]^{2}}$$

Explicit solution:

$$\hat{Z}_{S}(\underline{x}) = \underline{k}_{S}(\underline{x}) \cdot \mathbb{K}_{S}^{-1} \underline{Z}_{T}$$

where $\begin{cases} (\underline{k}_{S}(\underline{x}))_{\mu} = K_{S}(\underline{x}_{\mu}, \underline{x}) \\ & \text{kernel overlap} \\ (\mathbb{K}_{S})_{\mu\nu} = K_{S}(\underline{x}_{\mu}, \underline{x}_{\nu}) \\ & \text{Gram matrix} \end{cases}$

$$\frac{\hat{Z}_{S}(\underline{x}) = \sum_{\mu=1}^{n} c_{\mu} K_{S}(\underline{x}_{\mu}, \underline{x})}{\text{Minimize} = \frac{1}{n} \sum_{\mu=1}^{n} \left[\hat{Z}_{S}(\underline{x}_{\mu}) - Z_{T}(\underline{x}_{\mu}) \right]^{2}}$$

Explicit solution:

$$\hat{Z}_{S}(\underline{x}) = \underline{k}_{S}(\underline{x}) \cdot \mathbb{K}_{S}^{-1} \underline{Z}_{T}$$

where $\begin{cases} (\underline{k}_{S}(\underline{x}))_{\mu} = K_{S}(\underline{x}_{\mu}, \underline{x}) \\ & \text{kernel overlap} \end{cases} \\ (\mathbb{K}_{S})_{\mu\nu} = K_{S}(\underline{x}_{\mu}, \underline{x}_{\nu}) \\ & \text{Gram matrix} \end{cases} \\ (\underline{Z}_{T})_{\mu} = Z_{T}(\underline{x}_{\mu}) \\ & \text{training data} \end{cases}$

$$\frac{\hat{Z}_{S}(\underline{x}) = \sum_{\mu=1}^{n} c_{\mu} K_{S}(\underline{x}_{\mu}, \underline{x})}{\text{Minimize} = \frac{1}{n} \sum_{\mu=1}^{n} \left[\hat{Z}_{S}(\underline{x}_{\mu}) - Z_{T}(\underline{x}_{\mu}) \right]^{2}}$$

 $\begin{array}{l} \underline{\text{Explicit solution:}} \\ \hat{Z}_{S}(\underline{x}) = \underline{k}_{S}(\underline{x}) \cdot \mathbb{K}_{S}^{-1} \underline{Z}_{T} \end{array} \text{ where } \begin{cases} (\underline{k}_{S}(\underline{x}))_{\mu} = K_{S}(\underline{x}_{\mu}, \underline{x}) \\ (\mathbb{K}_{S})_{\mu\nu} = K_{S}(\underline{x}_{\mu}, \underline{x}_{\nu}) \\ \text{Gram matrix} \end{cases} \\ (\underline{Z}_{T})_{\mu} = Z_{T}(\underline{x}_{\mu}) \\ \text{training data} \end{cases}$

Compute the generalization error ϵ and how it scales with n

$$\epsilon = \mathbb{E}_{m{T}} \int \mathrm{d}^d \underline{x} \, \left[\hat{Z}_S(\underline{x}) - m{Z}_{m{T}}(\underline{x})
ight]^2 \sim n^{-eta}$$

14

To compute the generalization error:

- We look at the problem in the **frequency domain**
- We assume that $ilde{K}_S(\underline{w}) \sim \|\underline{w}\|^{-\alpha_S}$ and $ilde{K}_T(\underline{w}) \sim \|\underline{w}\|^{-\alpha_T}$ as $\|\underline{w}\| \to \infty$

To compute the generalization error:

- We look at the problem in the **frequency domain**
- We assume that $ilde{K}_S(\underline{w}) \sim \|\underline{w}\|^{-\alpha_S}$ and $ilde{K}_T(\underline{w}) \sim \|\underline{w}\|^{-\alpha_T}$ as $\|\underline{w}\| \to \infty$

E.g. Laplace has $\alpha = d + 1$ and Gaussian has $\alpha = \infty$

To compute the generalization error:

- We look at the problem in the **frequency domain**
- We assume that $ilde{K}_S(\underline{w}) \sim \|\underline{w}\|^{-\alpha_S}$ and $ilde{K}_T(\underline{w}) \sim \|\underline{w}\|^{-\alpha_T}$ as $\|\underline{w}\| \to \infty$

E.g. Laplace has $\alpha = d + 1$ and Gaussian has $\alpha = \infty$

• SIMPLIFYING ASSUMPTION: We take the *n* points \underline{x}_{μ} on a regular *d*-dim lattice!

To compute the generalization error:

- We look at the problem in the **frequency domain**
- We assume that $ilde{K}_S(\underline{w}) \sim \|\underline{w}\|^{-\alpha_S}$ and $ilde{K}_T(\underline{w}) \sim \|\underline{w}\|^{-\alpha_T}$ as $\|\underline{w}\| \to \infty$

E.g. Laplace has $\alpha = d + 1$ and Gaussian has $\alpha = \infty$

• SIMPLIFYING ASSUMPTION: We take the *n* points \underline{x}_{μ} on a regular *d*-dim lattice!

(details: arXiv:1905.10843)

Then we can show that

for $n\gg 1$ $\epsilon\sim n^{-eta}$ with $eta=rac{1}{d}\min(lpha_T-d,2lpha_S)$

$$eta = rac{1}{d}\min(lpha_T - d, 2lpha_S)$$

• Large $\alpha \rightarrow$ fast decay at high freq \rightarrow **indifference** to **local details**

• α_T is intrinsic to the **data** (T), α_S depends on the **algorithm** (S)

- If α_S is large enough, β takes the largest possible value $\frac{\alpha_T d}{d}$ (optimal learning)
- As soon as α_S is small enough, $\beta = \frac{2\alpha_S}{d}$

TEACHER-STUDENT: COMPARISON (1/2)

What is the prediction for our simulations? $eta = rac{1}{d}\min(lpha_T - d, 2lpha_S)$

• If Teacher=Student=Laplace

 $(\alpha_T = \alpha_S = d+1)$

$$eta=rac{lpha_T-d}{d}=rac{1}{d}$$
 (curse of dimensionality!)

• If Teacher=Gaussian, Student=Laplace $(\alpha_T = \infty, \alpha_S = d + 1)$

$$eta = rac{2lpha_S}{d} = 2 + rac{2}{d}$$

TEACHER-STUDENT: COMPARISON (2/2)

• Our result matches the numerical simulations

(on hypersphere)

• There are finite size effects (small *n*)



TEACHER-STUDENT: MATÉRN TEACHER

10²

n

MSE

10⁻⁸ -

 $v = 0.5, \ \beta = 1.0$

101

- $v = 1, \beta = 2$ - $v = 2, \beta = 4$ - $v = 4, \beta = 4$

Matérn kernels: $K_T(\underline{x}) = rac{2^{1u}}{\Gamma(
u)} z^
u \mathcal{K}_
u(z), \quad z = \sqrt{2
u} rac{\|\underline{x}\|}{\sigma}, \quad lpha = d + 2
u$ Laplace student, $K_S(\underline{x}) = \exp\left(-\frac{\|\underline{x}\|}{\sigma}\right)$ d = 110⁰ 10^{-2} $\beta = \min(2\nu, 4)$ 10^{-4} 10^{-6}

10³

NEAREST-NEIGHBOR DISTANCE

Same result with points on *regular lattice* or *random hypersphere*?

What matters is how **nearest-neighbor distance** δ scales with n (conjecture)



In both cases
$$\delta \sim n^{rac{1}{d}}$$

Finite size effects: asymptotic scaling only when n is large enough

What about real data?

 \longrightarrow second order approximation with a Gaussian process K_T :

does it capture some aspects?

What about real data?

 \longrightarrow second order approximation with a Gaussian process K_T : does it capture some aspects?

• Gaussian processes are *s*-times (mean-square) differentiable,

$$s=rac{lpha_T-d}{2}$$

What about real data? \longrightarrow second order approximation with a Gaussian process K_T : does it capture some aspects?

• Gaussian processes are *s*-times (mean-square) differentiable,

$$s=rac{lpha_T-d}{2}$$

• Fitted exponents are $\beta \approx 0.4$ (MNIST) and $\beta \approx 0.1$ (CIFAR10), regardless of the Student $\longrightarrow \beta = \frac{\alpha_T - d}{d}$

(since $\beta = rac{1}{d}\min(lpha_T - d, 2lpha_S)$ indep. of $lpha_S \longrightarrow eta = rac{lpha_T - d}{d}$)

What about real data? \longrightarrow second order approximation with a Gaussian process K_T : does it capture some aspects?

• Gaussian processes are *s*-times (mean-square) differentiable,

$$s=rac{lpha_T-d}{2}$$

• Fitted exponents are $\beta \approx 0.4$ (MNIST) and $\beta \approx 0.1$ (CIFAR10), regardless of the Student $\longrightarrow \beta = \frac{\alpha_T - d}{d}$

(since $\beta = rac{1}{d}\min(lpha_T - d, 2lpha_S)$ indep. of $lpha_S \longrightarrow eta = rac{lpha_T - d}{d}$)

 $\longrightarrow s = rac{1}{2}eta d$, s pprox 0.2d pprox 156 (MNIST) and s pprox 0.05d pprox 153 (CIFAR10)

This number is unreasonably large!

EFFECTIVE DIMENSION

• Measure NN-distance δ

ullet $\delta \sim n^{- ext{some exponent}}$



EFFECTIVE DIMENSION



EFFECTIVE DIMENSION



CURSE OF DIMENSIONALITY (1/2)

• Loosely speaking, the (optimal) exponent is

$$eta pprox rac{ ext{smoothness } lpha_T - d = 2s}{ ext{manifold dimension } d}$$

- To avoid the curse of dimensionality ($\beta \sim \frac{1}{d}$):
 - either the dimension of the manifold is small
 - or the data are extremely smooth

- Indeed, what happens if we consider a field $Z_T(\underline{x})$ that
 - is an instance of a Teacher K_T (α_T)
 - lies in the RKHS of a Student K_S (α_S)

- Indeed, what happens if we consider a field $Z_T(\underline{x})$ that
 - is an instance of a Teacher K_T (α_T)
 - lies in the RKHS of a Student K_S (α_S)

$$\begin{split} \mathbb{E}_{T} \| Z_{T} \|_{K_{S}} &= \\ \mathbb{E}_{T} \int \mathrm{d}^{d} \underline{x} \mathrm{d}^{d} \underline{y} \, Z_{T}(\underline{x}) K_{S}^{-1}(\underline{x}, \underline{y}) Z_{T}(\underline{y}) = \\ \int \mathrm{d}^{d} \underline{x} \mathrm{d}^{d} \underline{y} \, K_{T}(\underline{x}, \underline{y}) K_{S}^{-1}(\underline{x}, \underline{y}) < \infty \end{split} \qquad \Longrightarrow \qquad \alpha_{T} > \alpha_{S} + d \end{split}$$

- Indeed, what happens if we consider a field $Z_T(\underline{x})$ that
 - is an instance of a Teacher K_T (α_T)
 - lies in the RKHS of a Student K_S (α_S)

$$\begin{split} \mathbb{E}_{T} \| Z_{T} \|_{K_{S}} &= \\ \mathbb{E}_{T} \int \mathrm{d}^{d} \underline{x} \mathrm{d}^{d} \underline{y} Z_{T}(\underline{x}) K_{S}^{-1}(\underline{x}, \underline{y}) Z_{T}(\underline{y}) = \\ \int \mathrm{d}^{d} \underline{x} \mathrm{d}^{d} \underline{y} K_{T}(\underline{x}, \underline{y}) K_{S}^{-1}(\underline{x}, \underline{y}) < \infty \\ K_{S}(\underline{0}) \propto \int \mathrm{d} \underline{w} \, \tilde{K}_{S}(\underline{w}) < \infty \qquad \Longrightarrow \qquad \alpha_{S} > d \end{split}$$

- Indeed, what happens if we consider a field $Z_T(\underline{x})$ that
 - is an instance of a Teacher K_T (α_T)
 - lies in the RKHS of a Student K_S (α_S)

$$\begin{split} \mathbb{E}_{T} \| Z_{T} \|_{K_{S}} &= \\ \mathbb{E}_{T} \int \mathrm{d}^{d} \underline{x} \mathrm{d}^{d} \underline{y} \, Z_{T}(\underline{x}) K_{S}^{-1}(\underline{x}, \underline{y}) Z_{T}(\underline{y}) = \\ \int \mathrm{d}^{d} \underline{x} \mathrm{d}^{d} \underline{y} \, K_{T}(\underline{x}, \underline{y}) K_{S}^{-1}(\underline{x}, \underline{y}) < \infty \\ K_{S}(\underline{0}) \propto \int \mathrm{d} \underline{w} \, \tilde{K}_{S}(\underline{w}) < \infty \qquad \Longrightarrow \qquad \alpha_{S} > d \end{split}$$

(it scales with *d*!)

 $\longrightarrow \beta > \frac{1}{2}$

Therefore the smoothness must be $s = \frac{\alpha_T - d}{2} > \frac{d}{2}$

CURSE OF DIMENSIONALITY (2/2)

• Assume that the data are not smooth enough and live in *d* large

• **Dimensionality reduction** in the task rather than in the data?

• E.g. the n points \underline{x}_{μ} live in \mathbb{R}^d , but the target function is such that

$$Z_T(\underline{x}) = Z_T(\underline{x}_\parallel) \equiv Z_T(x_1,\ldots,x_{d_\parallel}), \quad d_\parallel < d$$

Similar setting studied in Bach 2017

• Can kernels understand the lower dimensional structure?

TASK INVARIANCE: KERNEL REGRESSION (1/2)

Theorem (informal formulation):

Similar result in Bach 2017

Two reasons contribute to this result:

- the nearest-neighbor distance always scales as $\delta \sim n^{-rac{1}{d}}$
- $\alpha_T(d) d$ only depends on the function $K_T(z)$ and not on d

TASK INVARIANCE: KERNEL REGRESSION (2/2)

Teacher = Matérn (with parameter ν), Student = Laplace, d=4



TASK INVARIANCE: CLASSIFICATION (1/2)

Classification with the margin SVM algorithm:

$$\hat{y}(\underline{x}) = ext{sign}\left[\sum_{\mu=1}^{n} c_{\mu} K\left(rac{\|\underline{x}-\underline{x}^{\mu}\|}{\sigma}
ight) + b
ight]$$

find $\{c_{\mu}\}, b$ by minimizing some function

We consider a very simple setting:

• the label is
$$y(\underline{x}) = y(x_1) \; \longrightarrow \; d_\parallel = 1$$





TASK INVARIANCE: CLASSIFICATION (2/2)

Vary kernel scale $\sigma \longrightarrow$ two regimes!

• $\sigma \ll \delta$: then the estimator is tantamount to a **nearest-neighbor algorithm** \longrightarrow curse of dimensionality $\beta = \frac{1}{d}$

• $\sigma \gg \delta$: important **correlations** in c_{μ} due to the **long-range kernel**. For the hyperplane with $d_{\parallel} = 1$ we find $\beta = O(d^0)!$

No curse of dimensionality!

THE NEAREST-NEIGHBOR LIMIT

hyperplane interface

using a Laplace kernel and varying the dimension d:

 $\beta = \frac{1}{d}$



KERNEL CORRELATIONS (1/2)

When $\sigma \gg \delta$ we can expand the kernel overlaps:

$$K\left(rac{\|\underline{x}-\underline{x}^{\mu}\|}{\sigma}
ight)pprox K(0)- ext{const} imes \left(rac{\|\underline{x}-\underline{x}^{\mu}\|}{\sigma}
ight)^{\xi}$$

(the exponent ξ is linked to the smoothness of the kernel)

We can derive some scaling arguments that lead to an exponent

$$eta = rac{d+\xi-1}{3d+\xi-3}$$

KERNEL CORRELATIONS (1/2)

When $\sigma \gg \delta$ we can expand the kernel overlaps:

$$K\left(rac{\|x-\underline{x}^{\mu}\|}{\sigma}
ight)pprox K(0)- ext{const} imes \left(rac{\|x-\underline{x}^{\mu}\|}{\sigma}
ight)^{\xi}$$

(the exponent ξ is linked to the smoothness of the kernel)

We can derive some scaling arguments that lead to an exponent

$$eta = rac{d+\xi-1}{3d+\xi-3}$$

Idea:

- support vectors ($c_{\mu}
 eq 0$) are close to the interface
- we impose that the decision boundary has $\mathcal{O}(1)$ spatial fluctuations on a scale proportional to δ

KERNEL CORRELATIONS (2/2)



KERN^{*n*}EL CORRELATIONS: HYPERSPHERE

What about other interfaces?

boundary = hypersphere:

$$egin{aligned} y(\underline{x}) &= ext{sign}(\|\underline{x}\|-R) \ (d_{\|} &= 1) \end{aligned}$$

$$eta = rac{d+\xi-1}{3d+\xi-3}$$

(same exponent!)

(similar scaling arguments apply, provided $R \gg \delta$)



CONCLUSION arXiv:1905.10843 + paper to be released soon!

• Learning curves of real data decay as **power laws** with exponents

 $\frac{1}{d} \ll \beta < \frac{1}{2}$

 We introduce a **new framework** that links the exponent β to the degree of smoothness of Gaussian random data

• We justify how different kernels can lead to the same exponent β

• We show that the **effective dimension** of real data is $\ll d$. It can be linked to a (small) **effective smoothness** *s*

 We show that kernel regression is not able to capture invariants in the task, while kernel classification can

(in some regime and for **smooth interfaces**)