

①

Lecture 6: Kernel methods

Last week: Lazy / linear regime

⇒ optimization regime where optimization is provably tractable

$$\hat{R}_m(\theta_t) \leq \hat{R}_m(\theta_0) e^{-c_0 t}$$

Success???

↳ what we actually care about: generalization (test) error $R(\theta_t)$

What is the performance of NNs trained in this regime?

(and in particular, can it explain the success of NNs in practice?)

In the lazy / linear regime, we saw that we can effectively replace the NN by its linearization around θ_0 :

$$f_{\text{lin}}(x; \theta) = f(x; \theta_0) + \langle \theta - \theta_0, \nabla_{\theta} f(x; \theta_0) \rangle$$

Denote $a = \theta - \theta_0$.

From lecture 2, GD converges to

$$\hat{a} = \operatorname{argmin} \left\{ \|a\|_2 : \langle a, \nabla_{\theta} f(x_i; \theta_0) \rangle = y_i - f(x_i; \theta_0) \right\}$$

(2)

and the model obtained is

$$\hat{f}(x) = f(x; \theta_0) + \langle \hat{a}, \nabla_{\theta} f(x; \theta_0) \rangle$$

We want to understand: What is the test error of \hat{f} ?

In fact, this is a special example of a "linear" or "kernel" method also called "kernel machines"

↳ methods that were developed in the 80 - 90's, which were state-of-the-art in ML before being replaced by NNs in the 2010s
(disappointing ???)

This connection between NNs trained by GD and kernel methods has spurred renewed interest in kernel methods

Are NNs simply glorified kernel methods?

Short answer:
no!!!

Lots of work since 2018 on understanding the fine grained performance of kernel methods. In particular

- interpolating solution
- non-monotonic behavior in learning curve

(both not covered by classical works on kernels)

3

Goal for today's lecture:

- * (Gentle) introduction to kernel methods
- * Sharp analysis of kernel ridge regression
- * Limitations of kernel methods

$$f_{\text{lin}}(x; \theta) = \underbrace{f(x; \theta_0)}_{\substack{\text{offset, not learned} \\ \text{random}}} + \underbrace{\langle \theta - \theta_0, \nabla_{\theta} f(x; \theta_0) \rangle}_{\text{"featureization of the data"}}$$

Linear model: $\Phi: \mathbb{R}^d \rightarrow \mathbb{R}^P$ } embedding data in a higher dimensional space
 $x \mapsto \Phi(x)$

Model: $f(x; \theta) = \langle \theta, \Phi(x) \rangle$ parameter $\theta \in \mathbb{R}^P$

↳ generalization of standard linear model $x \mapsto \langle \theta, x \rangle$

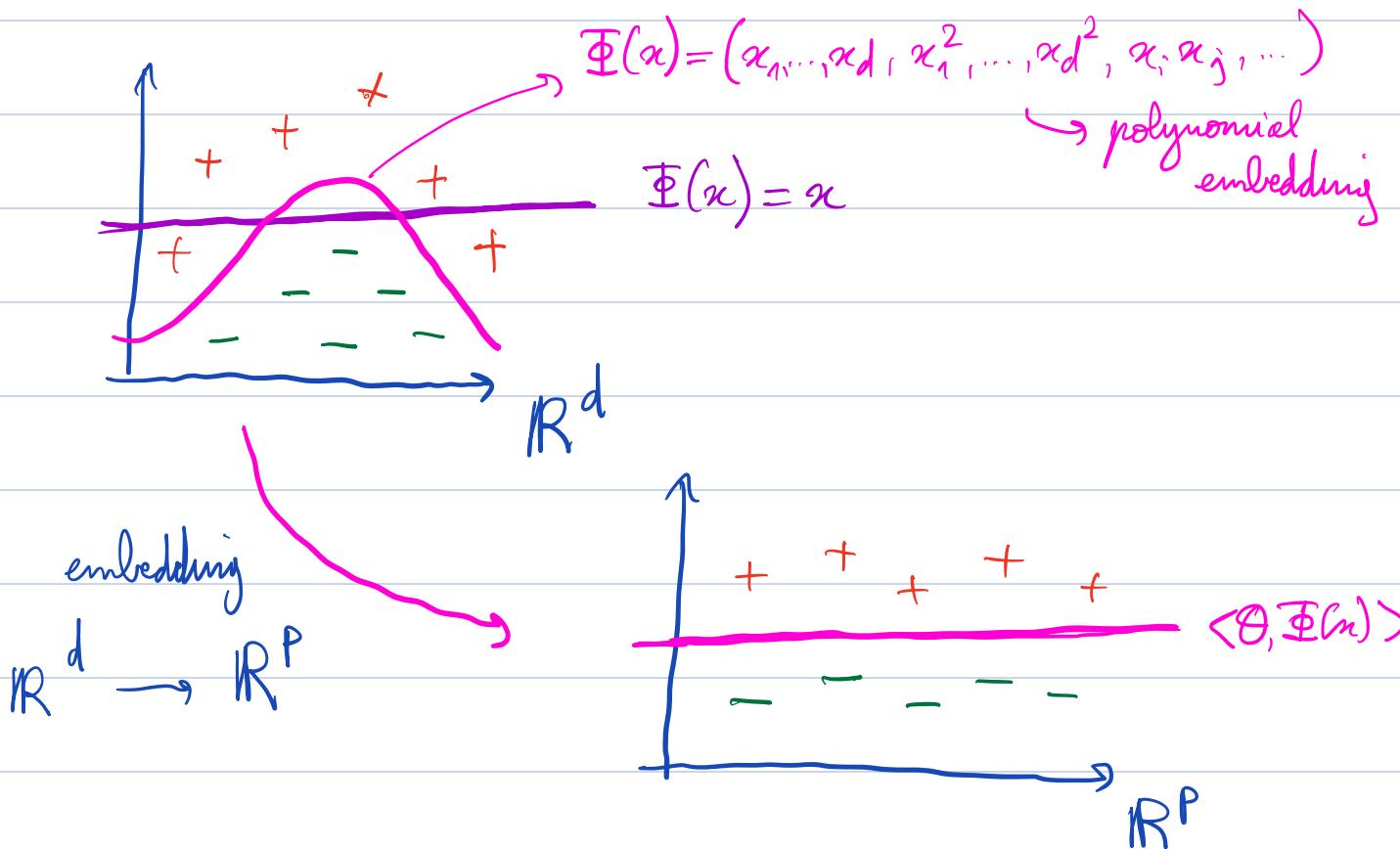
→ here still "linear model": linear in the parameters

but non linear in the data x

Why this is interesting?

4

E.g. (x_i, y_i) $y_i \in \{-1, 1\}$ predictor sign $\langle \theta, \Phi(x) \rangle$



→ can fit way more fits

→ in fact, typically $p=\infty$ and we can fit any function
(if we have enough data)

"universal"

(5)

Kernel methods

General definition:

* $(\mathcal{F}, \langle \cdot, \cdot \rangle_{\mathcal{F}})$ "feature space" with some inner product $\langle \cdot, \cdot \rangle_{\mathcal{F}}$
 norm $\|f\|_{\mathcal{F}} = \sqrt{\langle f, f \rangle_{\mathcal{F}}}$ (examples below)

* $\Phi: \mathcal{X} \rightarrow \mathcal{F}$ "featureization map" that embeds data x
 in $\Phi(x) \in \mathcal{F}$

* $\mathcal{H} = \{h_{\theta}: \mathcal{X} \rightarrow \mathbb{R} \quad h_{\theta}(x) = \langle \theta, \Phi(x) \rangle : \theta \in \mathcal{F}\}$

↳ space of functions which inherit scalar product/norm from \mathcal{F}

$$\langle h_{\theta}, h_{\theta'} \rangle_{\mathcal{H}} = \langle \theta, \theta' \rangle \quad \|h_{\theta}\|_{\mathcal{H}} = \|\theta\|_{\mathcal{F}}$$

Examples (1) $\mathcal{F} = \mathbb{R}^d$ with standard euclidean scalar product

$$\langle u, v \rangle = u_1 v_1 + \dots + u_d v_d$$

$$\Phi(x) = x$$

→ standard linear model $h_{\theta}(x) = \langle x, \theta \rangle \quad \|h_{\theta}\|_{\mathcal{H}} = \|\theta\|_2$

(6)

(2) $F = \mathbb{R}^{k+1}$ with standard euclidean scalar product

$$\Phi(x) = (1, x, x^2, \dots, x^k)$$

\Rightarrow standard polynomial regression model

$$h_\theta(x) = \theta_0 + \theta_1 x + \dots + \theta_{k+1} x^k$$

(3) (Ω, μ) probability space

$$L^2(\Omega, \mu) = \left\{ a: \Omega \rightarrow \mathbb{R}: \int a(\omega)^2 \mu(d\omega) < \infty \right\}$$

(space of squared integrable fcts on Ω)

$$\langle a, b \rangle_{L^2} = \int a(\omega) b(\omega) \mu(d\omega)$$

$$\|a\|_{L^2}^2 = \int a(\omega)^2 \mu(d\omega)$$

$$\Phi: \mathcal{X} \rightarrow L^2(\Omega, \mu) \quad x \mapsto \{ \omega \mapsto \sigma(x; \omega) \}$$

$$h_a(x) = \langle a, \Phi(x) \rangle_{L^2} = \int a(\omega) \sigma(x; \omega) \mu(d\omega)$$

$$\|h_a\|_{\mathcal{M}} = \|a\|_{L^2} \quad \mathcal{M} = \text{space of } \infty \text{ width 2-layer NNs with } \|a\|_{L^2} < \infty$$

(7)

Rmk: $(H, \langle \cdot, \cdot \rangle_H)$ forms a "Reproducing kernel Hilbert space" (RKHS)

which is a space uniquely defined by a positive semi-definite (PSD) kernel $K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

i.e.

$$\forall n \in \mathbb{N}, \forall x_1, \dots, x_n \in \mathcal{X}, \forall a_1, \dots, a_n \in \mathbb{R}$$

$$\sum_{i,j=1}^n a_i a_j K(x_i, x_j) \geq 0$$

In the above case: $K(x, x') = \langle \Phi(x), \Phi(x') \rangle_F$

$$(\text{e.g., } K(x, x') = \int \sigma(x; \omega) \sigma(x'; \omega) \mu(d\omega))$$

Those are two equivalent ways of defining $(H, \langle \cdot, \cdot \rangle_H)$

Given K , define $\Phi(x) = \{z \mapsto K(x, z)\} \in H$

$$K(x, x') = \langle \Phi(x), \Phi(x') \rangle_H \quad \text{"reproducing property"} \\ \langle h, \Phi(x) \rangle_H = h(x)$$

(I think defining RKHS through feature map

is the most illuminating. Δ some technical difficulties
that I skipped

Δ Given H or K : there are infinite number of possible feature maps which yield same H , but not necessarily equivalent when doing approximation.

(8)

Kernel methods

$$\hat{\theta} = \underset{\theta \in \mathcal{F}}{\operatorname{argmin}} \left\{ \frac{1}{m} \sum_{i=1}^m l(y_i, \langle \theta, \phi(x_i) \rangle_F) + \gamma \|\theta\|_F^2 \right\}$$

↔

$$\hat{h} = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \left\{ \frac{1}{m} \sum_{i=1}^m l(y_i, h(x_i)) + \gamma \|h\|_H^2 \right\}$$

Is it tractable? $\Phi(x_i) \rightarrow$ potentially infinite dimensional
 optimum θ : ∞ -dimensional problem

[doesn't matter if \mathcal{H} is way more expressive than
 standard linear models if we can't train them efficiently]

"Kernel trick": in fact, reduces to n -dimensional convex problem!

↳ This is the reason kernel methods become so popular

Thm [Representer Theorem] For any loss l (doesn't need to be convex!)

solut° $\hat{\theta}$ can be written as

$$\hat{\theta} = \sum_{i=1}^n a_i \Phi(x_i)$$

Rank: To find $\hat{\Theta}$, it suffices to restrict $\Theta \in \text{span}\{\Phi(x_i) : i \in [n]\}$ (9)

\hookrightarrow n -dim linear subspace

Recall $K(x, x') = \langle \Phi(x), \Phi(x') \rangle_F$

Denote $k_m(x) = (K(x, x_1), \dots, K(x, x_m)) \in \mathbb{R}^m$

$K_m = (K(x_i, x_j))_{i,j=1}^m$ "Empirical kernel matrix"

Then

$$\hat{a} = \underset{a \in \mathbb{R}^m}{\operatorname{argmin}} \left\{ \frac{1}{m} \sum_{i=1}^m l(y_i, a^T k_m(x_i)) + \lambda a^T K a \right\}$$

\rightarrow If l convex then this is a convex optimization problem over m variables

Solution $\hat{\Theta} = \sum_{i=1}^m \hat{a}_i \Phi(x_i)$

$$\hat{h}(x) = \langle \hat{\Theta}, \Phi(x) \rangle = \sum_{i=1}^m \hat{a}_i K(x, x_i)$$

Hence: for training and prediction, we never need to compute $\Phi(x)$, just the inner-product $K(x, x') = \langle \Phi(x), \Phi(x') \rangle_F$ which is often easy!

This is known as the "kernel trick".

(10)

Proof: Denote $V = \text{span} \{ \Phi(x_i) : i \in [m] \}$ linear subspace in \mathcal{F}

$$\mathcal{F} = V \oplus V_{\perp} \quad V_{\perp} = \{ \theta \in \mathcal{F} : \langle \theta, v \rangle_{\mathcal{F}} = 0 \text{ for all } v \in V \}$$

= orthogonal space

$$\theta = \theta_{\parallel} + \theta_{\perp} \quad \|\theta\|_{\mathcal{F}}^2 = \|\theta_{\parallel}\|_2^2 + \|\theta_{\perp}\|_2^2$$

$$\min_{\theta_{\parallel}, \theta_{\perp}} \frac{1}{m} \sum_{i=1}^m l(y_i, \langle \theta_{\parallel} + \theta_{\perp}, \phi(x_i) \rangle) + \lambda \|\theta_{\parallel}\|_{\mathcal{F}} + \lambda \|\theta_{\perp}\|_{\mathcal{F}}$$

θ_{\perp} only appears here so $\theta_{\perp} = 0$ to minimize

□

Kernel Ridge Regression

A popular example is KRR: $l(y, \hat{y}) = (y - \hat{y})^2$

$$\hat{h} = \underset{h \in H}{\operatorname{argmin}} \left\{ \sum_{i=1}^m (y_i - h(x_i))^2 + \lambda \|h\|_H^2 \right\}$$

$$\hat{h} = \sum_{i=1}^m \hat{\alpha}_i K(x, x_i) \quad k_m(x) = (K(x, x_i))_{i=1}^m$$

$$K_m = (K(x_i, x_j))_{i,j=1}^m$$

$$\|h\|_H^2 = \alpha^T K_m \alpha \quad y = (y_1, \dots, y_m)$$

$$\hat{\alpha} = \underset{\alpha \in \mathbb{R}^m}{\operatorname{argmin}} \left\{ \|y - K_m \alpha\|_2^2 + \lambda \alpha^T K_m \alpha \right\}$$

$$= (K_m^2 + \lambda K_m)^{-1} K_m y \stackrel{\uparrow}{=} (K_m + \lambda I_m)^{-1} y$$

KKT condition
(assume K_m is full rank)

$$\hat{h}(x) = k_m(x)^T (K_m + \lambda I_m)^{-1} y$$

Rmk: [“linear regression on Hilbert space”] $\hat{h}(x) = \langle \hat{\theta}, \Phi(x) \rangle$

$$\hat{\theta} = \Phi(X)^T (\Phi(X) \Phi(X)^T + \lambda I_m)^{-1} y$$

$$\Phi(X) = \begin{bmatrix} -\phi(x_1) - \\ \vdots \\ -\phi(x_m) - \end{bmatrix} \in \mathbb{R}^{n \times p}$$

Rmk: [Interpolating solution] $\downarrow \downarrow 0$

$$\hat{h} = \text{argmin} \left\{ \|h\|_H : h(x_i) = y_i \quad i \leq m \right\}$$

$$\hat{h}(x) = k_m(x) K_m^{-1} y$$

Diagonalization of the kernel

The performance of kernel methods depend crucially on the eigendecomposition of the kernel

$$\alpha \sim (\mathcal{X}, \nu) \quad L^2(\mu) = \left\{ f: \mathcal{X} \rightarrow \mathbb{R}, \|f\|_{L^2}^2 = \int f(x)^2 \nu(dx) < \infty \right\}$$

Can associate to K a linear operator $\mathbb{K}: L^2(\nu) \rightarrow L^2(\nu)$

$$(\mathbb{K}f)(x) = \int K(x, x') f(x') \nu(dx')$$

Assume $E_\alpha [K(\alpha, \alpha)] < \infty$ (K is "kern class")

Then by spectral theorem over bounded linear operator

$$\mathbb{K} = \sum_{j=1}^{\infty} \lambda_j \phi_j \phi_j^*$$

$\lambda_1 \geq \lambda_2 \geq \dots$ are the eigenvalues (non-negative)

ϕ_1, ϕ_2, \dots are the eigenfunctions

$$\mathbb{K} \phi_j = \lambda_j \phi_j$$

$$\langle \phi_j, \phi_k \rangle_{L^2(\nu)} = \delta_{jk}$$

In terms of the kernel function, this implies

$$K(x, x') = \sum_{j=1}^{\infty} \lambda_j \phi_j(x) \phi_j(x')$$

This suggests a "canonical" way of constructing an embedding:

$$\Phi(x) = \left(\lambda_j^{1/2} \phi_j(x) \right)_{j=1}^{\infty}$$

$\Phi(x) \in (l_2, <, >_{l_2})$ the Hilbert space of squared summable sequences

$$l_2 = \left\{ (\alpha_1, \alpha_2, \dots) : \sum_{i=1}^{\infty} \alpha_i^2 < \infty \right\}$$

$$\langle a, b \rangle_{l_2} = \sum_{i=1}^{\infty} a_i b_i$$

$$H = \left\{ h_{\theta}(x) = \langle \theta, \Phi(x) \rangle = \sum_{j=1}^{\infty} \theta_j \lambda_j^{1/2} \phi_j(x) : \theta \in l_2 \right\}$$

Rmk: If $\{\phi_j\}_{j=1}^{\infty}$ is a basis of $L^2(\nu)$, then for all

$$f \in L^2(\nu), \quad f = \sum_{j=1}^{\infty} \beta_j \phi_j = \sum_{j=1}^{\infty} \frac{\beta_j}{\lambda_j^{1/2}} \lambda_j^{1/2} \phi_j \xrightarrow{\langle \theta, \Phi(x) \rangle}$$

$$\|f\|_{L^2}^2 = \sum_{j=1}^{\infty} \beta_j^2$$

$$\|f\|_H^2 = \sum_{j=1}^{\infty} \frac{\beta_j^2}{\lambda_j} = \|\theta\|_{l_2}^2$$

Test error :

Bias - Variance decomposition

$$(x_i, y_i)_{i=1}^n \quad y_i = h_*(x_i) + \varepsilon_i \quad \varepsilon_i \sim N(0, \sigma_\varepsilon^2)$$

$$\begin{aligned} \text{(excess test error)} \quad R(\hat{h}) &= \mathbb{E}_x \left[(h_*(x) - \hat{h}(x))^2 \right] \\ &= \mathbb{E}_x \left[(h_*(x) - k_m(x)^T (K_m + \lambda)^{-1} y)^2 \right] \end{aligned}$$

$$\mathbb{E}_\varepsilon [R(\hat{h})] = B_m(\lambda) + V_m(\lambda)$$

$$V_m(\lambda) = \sigma_\varepsilon^2 \mathbb{E}_x \left[k_m(x) k_m(x)^T \right]$$

$$B_m(\lambda) = \langle h_{*,m}, (K_m + \lambda)^{-1} M (K_m + \lambda)^{-1} h_{*,m} \rangle$$

$$- 2 \langle v, (K_m + \lambda)^{-1} h_{*,m} \rangle + \|h_*\|_{L^2}^2$$

where $M_{ij} = \mathbb{E}_x [K(x_i, x_j) K(x_i, x_j)]$

$$v_i = \mathbb{E}_x [K(x_i, x) h_*(x)]$$

$$h_{*,m} = (h_*(x_1), \dots, h_*(x_m))$$

Do
exercise
at
home !!

Sharp characterization of Bias / Variance

Let us use the "canonical factorization" from eigendecomposition

$$\Phi(x) = \left(\lambda_j^{\frac{1}{2}} \phi_j(x) \right)_{j=1}^{\infty} \quad y_i = \langle \theta_*, \Phi(x_i) \rangle + \varepsilon_i$$

$$z_i := \Phi(x_i) \quad Z = \begin{bmatrix} z_1 \\ \vdots \\ z_m \end{bmatrix} \in \mathbb{R}^{m \times \infty}$$

$$z_i \text{ iid} \quad \mathbb{E}[zz^T] = \Sigma = \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \end{pmatrix}$$

$$B_n(\lambda) = \lambda^2 \langle \theta_*, (Z^T Z + \lambda)^{-1} \Sigma (Z^T Z + \lambda)^{-1} \theta_* \rangle$$

$$V_n(\lambda) = \sigma_\varepsilon^2 \text{Tr}(\Sigma Z^T Z (Z^T Z + \lambda)^{-2})$$

Using random matrix theory, we have very sharp characterization of these quantities

→ well approximated by "deterministic equivalents"

Define "effective regularization" = λ_* as unique solution of fixed point equation:

$$m - \frac{1}{\lambda_*} = T(\Sigma (\Sigma + \lambda_*)^{-1})$$

Define:

$$\bar{V}_m(\lambda_*) = \sigma_\varepsilon^2 \cdot \frac{T(\Sigma^2 (\Sigma + \lambda_*)^{-2})}{m - T(\Sigma^2 (\Sigma + \lambda_*)^{-2})}$$

$$\bar{B}_m(\lambda_*) = \frac{\lambda_*^2 \langle \theta_*, \Sigma (\Sigma + \lambda_*)^{-2} \theta_* \rangle}{1 - \frac{1}{m} T(\Sigma^2 (\Sigma + \lambda_*)^{-2})}$$

$$\bar{R}_m(\lambda) = \bar{B}_m(\lambda_*) + \bar{V}_m(\lambda_*)$$

\hookrightarrow deterministic quantity that only depend on m, λ, Σ and θ_* (coefficients of the target fct)

Can show: $B_m(\lambda) = (1 + o_m(1)) \bar{B}_m(\lambda_*)$

$$V_m(\lambda) = (1 + o_m(1)) \bar{V}_m(\lambda_*)$$

Assumption [Hanson-Wright] Assume that for every P.S.D matrix $A \in \mathbb{R}^{n \times n}$ s.t $\text{Tr}(\Sigma A) < \infty$

$$P(|y^T A y - \text{Tr}(A\Sigma)| \geq t \cdot \|\Sigma^{\frac{1}{2}} A \Sigma^{\frac{1}{2}}\|_F) \leq C e^{-ct}$$

E.g. $\Sigma^{-\frac{1}{2}} y$ is iid subGaussian or has log-concave density
 ↳ can relax this assumption

Thm [Cheng, Montanari, '24, M., Saeed, '24] (sketch)

Under above assumption, we have with high prob.,

$$|B_m(\lambda) - \bar{B}_m(\lambda_*)| \leq \frac{C}{\sqrt{m}} \bar{B}_m(\lambda_*)$$

$$|V_m(\lambda) - \bar{V}_m(\lambda_*)| \leq \frac{C}{m} \bar{V}_m(\lambda_*)$$

Rmk: 1) Relative error: $R_m(\lambda) = \bar{R}_m(\lambda)(1 + o_m(1))$
 ↳ very sharp prediction and extremely good agreement with numerical simulations

2) Dimension-free: applies to $p=\infty$

Rank: [Gaussian sequence model] Equivalence with
following model:

Observe $y_j^s = \gamma_j^{\frac{1}{2}} \theta_{*,j} + \frac{\omega}{\sqrt{n}} g_j \quad j=1,2,\dots$

$$y^s = \sum \gamma_j^{\frac{1}{2}} \theta_{*,j} + \frac{\omega}{\sqrt{n}} g \quad g_j \sim N(0, I_p)$$

Goal: recover θ_* from y^s

$$\hat{\theta}^s = \underset{\theta \in R^\infty}{\operatorname{argmin}} \left\{ \|y^s - \sum \gamma_j^{\frac{1}{2}} \theta_j\|_2^2 + \gamma_* \|\theta\|_2^2 \right\}$$

ω fixed pt of: $\omega^2 = \sigma_\varepsilon^2 + E_g [\|\hat{\theta}^s - \theta_*\|_2^2]$

$$R_m^s = E_g [\|\hat{\theta}^s - \theta_*\|_2^2] = \bar{B}_m(\gamma_*) + \bar{V}_m(\gamma_*)$$

↑ same as KRR

	Original problem	Sequence model
design	X random	$\Sigma^{\frac{1}{2}}$ deterministic
reg.	λ	$\lambda_*(\lambda)$
noise	σ_ε^2	ω^2

Rmk: Even when $\lambda = 0^+$ (interpolating solution)

$\lambda_*(0^+) > 0$: self induced regularization

Model has more regularization than what we naively expect from $\lambda = 0$

Rmk: Test error \approx det equiv : doesn't depend on particular data distribution, only on eigenvalues

→ "Gaussian universality" or just "universality"

Inner-Product kernel

Last week: 2 layer neural network

e.g. $f_{RF}(x; \alpha) = \langle \alpha, \Phi_{RF}(x) \rangle$

$$\Phi_{RF}(x) = \frac{1}{\sqrt{M}} \begin{pmatrix} \sigma(\langle x, \omega_1^0 \rangle) \\ \vdots \\ \sigma(\langle x, \omega_M^0 \rangle) \end{pmatrix}$$

Linearizing 2nd
layer weights

Associated kernel: $\omega_0 \sim N(0, \text{Id}_d)$

$$K_M(x_1, x_2) = \frac{1}{M} \sum_{j \in [M]} \sigma(\langle \omega_j^0, x_1 \rangle) \sigma(\langle \omega_j^0, x_2 \rangle)$$

$$\xrightarrow[M \rightarrow \infty]{} \mathbb{E}_{\omega^0 \sim N(0, \text{Id}_d)} [\sigma(\langle \omega, x_1 \rangle) \sigma(\langle \omega, x_2 \rangle)]$$

$$\|x_1\|_2 = \|x_2\|_2 = 1 \quad \tau = \langle x_1, x_2 \rangle$$

$$= \mathbb{E}_{G_1, G_2 \sim N(0, 1)} [\sigma(G_1) \sigma(t G_1 + \sqrt{1-t^2} G_2)]$$

$$= h(\langle x_1, x_2 \rangle)$$

More generally limiting kernel as $M \rightarrow \infty$ in fully connected NNs with Gaussian initialization are inner product kernels

What is generalization error of inner-product kernel?

$$(x_i, y_i)_{i=1}^m \quad x_i \sim \text{Unif}(\mathbb{S}^{d-1})$$

$$y_i = f_*(x_i) + \varepsilon_i \quad f_* \in L^2$$

We saw that the test error depends crucially on the eigendecomposition of the kernel.

For data uniformly distributed on the sphere, there is an explicit eigendecomposition in terms of spherical harmonics

$$L^2(\mathbb{S}^{d-1}) = \bigoplus_{k=0}^{\infty} V_k \quad \begin{array}{l} \rightarrow \text{space of } d^k \text{ spherical harmonics} \\ \text{i.e. space of polynomials of degree } k \\ \text{perpendicular to all polynomials of degree } \leq k-1. \end{array}$$

$$B_{d,k} := \dim(V_k) = \frac{d+2k-2}{k} \binom{d+k-3}{k-1} = \Theta(d^k)$$

→ basis $\{Y_{ks} : 1 \leq s \leq B_{d,k}\}$

then

$$h(\langle x_1, x_2 \rangle) = \sum_{h=0}^{\infty} \mathfrak{Z}_h \sum_{s=1}^{B_{d,k}} Y_{hs}(x_1) Y_{hs}(x_2)$$

each eigenvalue \mathfrak{Z}_h has degeneracy $B_{d,k}$

→ eigenspace associated to \mathfrak{Z}_h is V_h

$$T_h(K) = \mathbb{E}_{\alpha}[K(x, \alpha)] = h(1) = \sum_{h=0}^{\infty} \mathfrak{Z}_h B_{d,k}$$

meaning that $\mathfrak{Z}_h \approx \frac{1}{B_{d,k}} \approx d^{-h}$

We are interested in the test error in the polynomial
high-dimensional regime where $m, d \rightarrow \infty$

$$\frac{m}{d^k} \rightarrow \gamma \quad k, \gamma > 0$$

Consider empirical kernel matrix $K_m = (K(x_i, x_j))_{i,j=1}^m$

$$K_m = \sum_{k=0}^{\infty} \beta_k Y_k Y_k^T$$

where

$$Y_k = \begin{bmatrix} Y_{k,1}(x_1) & \dots & Y_{k,B_{d,k}}(x_1) \\ Y_{k,1}(x_m) & \dots & Y_{k,B_{d,k}}(x_m) \end{bmatrix} \in \mathbb{R}^{m \times B_{d,k}}$$

Case 1: $l < k < l+1$

$$K_m = \underbrace{\sum_{k=0}^l \beta_k Y_k Y_k^T}_{\text{K}_m \leq l} + \underbrace{\sum_{k=l+1}^{\infty} \beta_k Y_k Y_k^T}_{\text{K}_m > l}$$

* $\boxed{k \leq l}$: "low degree part" $\beta_k m \gtrsim d^{k-l} \rightarrow \infty$ as $d \rightarrow \infty$

$$\lambda_{\min}(Y_k Y_k^T) \asymp m \quad \beta_k m \gtrsim d^{k-l} \rightarrow \infty \text{ as } d \rightarrow \infty$$

hence eigenvalues of this part $\rightarrow \infty$

furthermore rank $1 + B_{d,1} + B_{d,2} + \dots + B_{d,l} = O(d^l)$

$\ll m$

Mence $K_{m,\leq l}$ low rank spiked metric

→ Thanks to that, KRR will fit perfectly
 $P_{\leq l} f_*$ the projection of f_* on degree- l polynomials

(i.e., best degree- l polynomial approx of f_*)

$$* \boxed{k \geq l+1} \quad Y_k \in \mathbb{R}^{m \times d^k} \quad \frac{Y_k Y_k^T}{B_{d,k}} \approx I_m$$

$$K_{m,>l} \approx \left(\sum_{h=l+1}^{\infty} \underbrace{\zeta_h}_{\mu_*} B_{d,h} \right) I_m$$

μ_* "self-induced regularization"

→ Because of that, KRR will not fit at all
 the higher degree part of f_*

$$\boxed{\text{Mence}} \quad m = d^k \quad l < k < l+1$$

$$\hat{f}_{\text{KRR}}(x) = P_{\leq l} f_*(x) + \underbrace{\Delta(x)}_{L^2 \text{ vanishing}}$$

→ can interpolate data if $\Delta = 0^+$

$$R(\hat{f}_{KRR}) = \|P_{\geq l} f_*\|_{L^2}^2 + o_d(1)$$

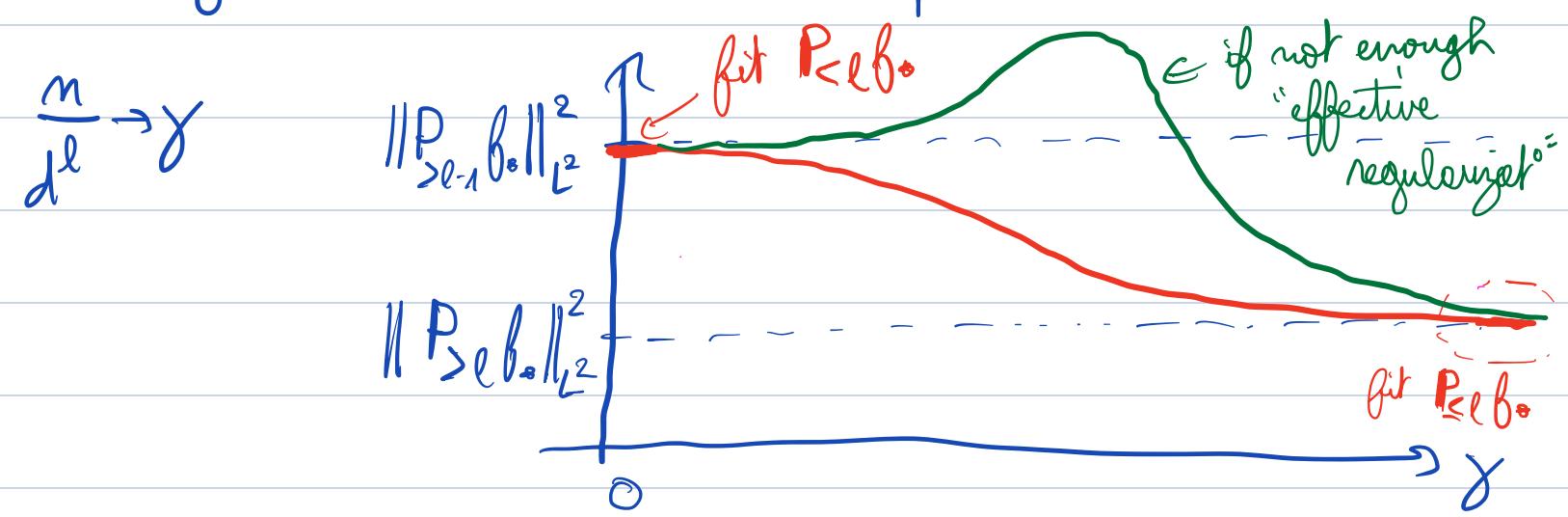
Case 2: $m \leq d^l$

$$K_m = \underbrace{K_{m, \leq l-1}}_{\substack{\text{low rank} \\ \rightarrow \text{fit perfectly}}} + \underbrace{\beta_l Y_l Y_l^T}_{\substack{\text{spikes} \\ P_{\leq l-1} f_*}} + \underbrace{K_{m, >l}}_{\substack{\approx \mu_0 I_m \\ \text{does not} \\ \text{fit } P_{\geq l} f_* \\ \text{at all}}}$$

$\beta_l \frac{Y_l Y_l^T}{B_{d,l}}$ behave as a Wishart matrix

$\approx (1)$

using RMT we can compute the test error



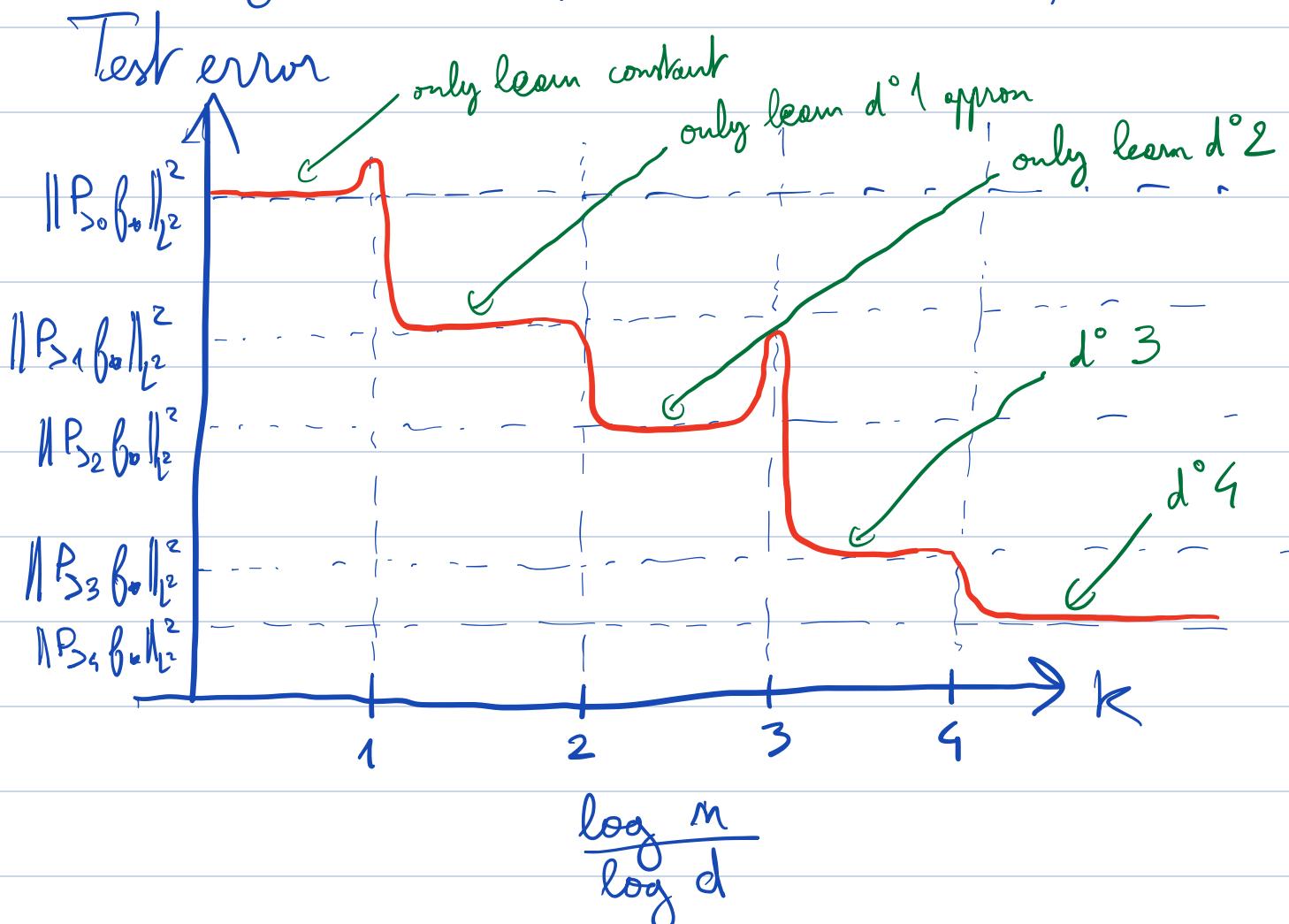
(27)

$$\frac{m}{d} \rightarrow \gamma$$

$\gamma \downarrow 0$ doesn't fit $P_{\ell} f_*$ at all

$\gamma \rightarrow \infty$ fit $P_{\ell} f_*$ perfectly

In between, if not enough "effective regularization"
there might be a peak at $m = B_{d,l}$



[Ghorbani et al '21]

[Misiakiewicz, '22]

Conclusion: let's say $d = 1000$ (typical in applications)

Then to get $d^{\circ 1}$ approx $n \approx 1000$

$d^{\circ 2}$ approx $n \approx 1\,000\,000$

$d^{\circ 3}$ approx $n \approx 1\,000\,000\,000$

Quite bad!!! Curse of dimensionality !!!



Often data is not isotropic but is concentrated on a smaller subspace of dimension d_{eff}

→ in that case $n \approx d_{\text{eff}}^k$ for $d^{\circ k}$ approx

E.g. CIFAR-10 $d_{\text{eff}} \approx 30$
 $n = 60\,000$

→ can only get $d^{\circ 3}$ approx

Power and limitations of kernel methods

What KRR does high level:

$$K = \sum_{j=1}^{\infty} \lambda_j \phi_j \phi_j^* \quad \lambda_1 \geq \lambda_2 \geq \dots$$

Target fact: $f_* = \sum_{j=1}^{\infty} \beta_j \phi_j$ and on data pts

$$\Rightarrow \text{KRR } \hat{f} \approx \sum_{j=1}^m \beta_j \phi_j$$

i.e. fit the project of f_* on top on eigenspaces

E.g. To fit d^k polynomial approximat: subspace of $\dim \mathcal{H}(d^k)$
 ↳ need $m \gtrsim d^k$
 → KRR with inner product kernel

This is minimax optimal: to fit all d^k polynomial we need
 $m = \sqrt{2}(d^k)$
 ↳ no method can do better!

Classical results on rates of kernels

$$(1) \mathcal{E}_L = \left\{ f_\infty : \mathbb{R}^d \rightarrow \mathbb{R} \text{ L-lipschitz} \right\}$$

$$\sup_{f_\infty \in \mathcal{E}_L} R_{\text{test}}(f_\infty, \hat{f}) \asymp n^{-1/d}$$

↳ worst-case, to get error $\leq \varepsilon$ need $n \geq \left(\frac{1}{\varepsilon}\right)^d$

CURSE OF DIMENSIONALITY

→ no methods can do better (minimax optimal)

$$(2) \mathcal{E}_S = \left\{ f_\infty \text{ with } S \text{ first derivatives bounded} \right\}$$

$$\sup_{f_\infty \in \mathcal{E}_S} R_{\text{test}}(f_\infty, \hat{f}) \asymp n^{-S/(S+d)}$$

$$\text{To get } \varepsilon \text{ error } n \asymp \left(\frac{1}{\varepsilon}\right)^{1+\frac{d}{S}}$$

if $S \asymp d$

↓
no curse of dim

but needs to be extremely smooth

→ again no method can do better

↳ KRR minimizes optimal

Hence:

KRR is adaptive to smoothness of the fct
 ↳ smoother fct will be easier to fit

E.g.: • $S = d$ smooth $n \asymp \left(\frac{1}{\epsilon}\right)^C$

• d^l polynomial $n \asymp d^l$

KRR adaptive to other properties of target fct?

E.g. $n \geq d^l$ to fit any degree- l polynomial no matter the polynomial!!

Other interesting classes of fcts: e.g. $f(x) = g(U^T x)$
 ("Multi-index fcts") $U \in \mathbb{R}^{d \times p}$

i.e. fcts that only depend on a low dimensional project of the data

minimax rate is $n^{-\Theta(\frac{1}{p})}$

→ KRR will not reach this minimax rate

(e.g. inner-product kernel doesn't care at all that only depend on low dim project)

Kernel methods will not be adaptive to low-dim structure

At the same time for NNs: $f_{NN}(x; \theta) = \sum a_j \sigma(\langle w_j, x \rangle)$

Idea: can align w_j with support \cup , then the model effectively dimension p and we can achieve minimax rate

$$\hat{f}_{\text{ERM}} = \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (y_i - f_{NN}(x_i; \theta))^2 + \lambda \|f\|_{F_1}^2$$

then $R(f_*, \hat{f}_{\text{ERM}}) \asymp n^{-\Theta(\frac{1}{p})}$

Of course, we would need to show this rate for a NN trained by SGD!

⇒ See next two lectures!!!

Summary

* Lazy regime :

→ show tractability of optimization ✓

→ NNs in this regime collapse on kernel methods

* Kernel methods are (relatively) well understood learning methods that have limited adaptivity

↳ heuristically, from n samples, can only fit fits in
a n -dimensional subspace $\text{span}\{\Phi(x_i)\}_{i=1}^n$ (indep of labels)

↳ "fixed features methods"

$$f(x) = \langle \theta, \underbrace{\Phi_{\theta^0}(x)}_{\text{NT model at init}} \rangle$$

→ NT model at init

"Fixed embedding" =

"Fixed representation" =

→ suffer curse of dimensionality on interesting feature classes
that can be approximated efficiently by NNs

* When NNs are trained beyond lazy regime
 w^t moves away from w^* .

$$f_t(x) = \langle \theta, \underline{\Phi_{w^t}(x)} \rangle$$

↳ embedding changes

→ NNs can learn a good representation of the data

beyond lazy regime : "feature learning" or "rich" requires

E.g. $a^T \sigma(W_t \alpha)$

W_t can align with the support of target fct $g(V^T x)$

↳ then fitting α becomes much easier

(effectively p -dimensional problem instead of)
 d -dimensional problem

In the next 2 lectures, we will show such an example of optimization regime with feature learning.

Dimension lower bound for Kernel methods

Below, I give a very simple lower bound on learning with kernel methods, which I think encapsulates the limitation of learning with kernels

- we can only fit a subspace of dimension n independent of the target function
- no adaptivity

Prop [Mru, '22, Abbe, Boix-Adserà, Mirekiewicz '22]

Let \mathcal{R} a Hilbert space with inner product $\langle \cdot, \cdot \rangle_{\mathcal{R}}$

Let $\{f_1, \dots, f_M\} \subset \mathcal{R}$ with $\|f_i\|_{\mathcal{R}} = 1$

Let $T \subset \mathcal{R}$ linear subspace of \mathcal{R} of dimension n

$$\text{Let } \varepsilon := \frac{1}{M} \sum_{i \in [M]} \inf_{g \in T} \|g - f_i\|_{\mathcal{R}}^2 \quad) \begin{array}{l} \text{average} \\ \text{approximation} \\ \text{error} \\ \text{of } f_i \text{ by } T \end{array}$$

Then

$$n \geq \frac{M}{\|G\|_{op}} (1 - \varepsilon)$$

$$\text{where } G = (\langle \beta_i, \beta_j \rangle_{\mathcal{R}})_{i,j=1}^M$$

Rank: What does it mean?

Consider $\mathcal{R} = L^2(\mathcal{X})$

When fitting with a kernel method with n data points

$$\hat{f} = \sum_{i=1}^n a_i K(x, x_i)$$

Take $T = \text{span} \{ x \mapsto K(x, x_i) : i \in [n] \}$

$$\hookrightarrow \dim(T) = n$$

Test error when fitting target fct $f_i \geq \inf_{g \in T} \|g - f_i\|_{L^2}^2$

Hence if test error $\leq \varepsilon$ for M orthogonal fcts f_i
we must have

$$n \geq M(1-\varepsilon)$$

Proof: Let ϕ_1, \dots, ϕ_n be an orthonormal basis of T and let Π_T be the orthogonal projection onto T in $(\mathbb{R}, \langle \cdot, \cdot \rangle_R)$.

In particular $\Pi_T f = \sum_{j=1}^n \phi_j \langle \phi_j, f \rangle_R$

$$\varepsilon = \frac{1}{M} \sum_{i=1}^M \inf_{g \in T} \|g - f_i\|_R^2 \quad (\text{definition})$$

$$= \frac{1}{M} \sum_{i=1}^M 1 - \|\Pi_T f_i\|_R^2 \quad (\|f_i\|_R^2 = 1)$$

$$= 1 - \frac{1}{M} \sum_{i=1}^M \sum_{j=1}^n \langle \phi_j, f_i \rangle_R^2$$

Use that

$$\begin{aligned} \sum_{i=1}^M \langle g, f_i \rangle_R^2 &= \left\langle g, \sum_{i=1}^M \langle g, f_i \rangle f_i \right\rangle_R \\ &\leq \|g\|_R \left(\sum_{i,j=1}^M \langle g, f_i \rangle \langle g, f_j \rangle \langle f_i, f_j \rangle \right)^{1/2} \\ &= \|g\|_R (b^T G b)^{1/2} \leq \|g\|_R \|G\|_{op}^{1/2} \|b\|_2 \end{aligned}$$

where $b = (\langle g, f_i \rangle)_{i=1}^M$

Note that $\|b\|_2^2$ is the left-hand side of
the above inequality, hence

(38)

$$\sum_{i=1}^M \langle g, f_i \rangle_R^2 \leq \|g\|_R^2 \|G\|_{op}$$

$$\text{Hence } 1 - \varepsilon = \frac{1}{M} \sum_{j=1}^n \sum_{i=1}^M \langle \phi_j, f_i \rangle^2 \leq \frac{n}{M} \|G\|_{op}$$

$$\Rightarrow n \geq \frac{M}{\|G\|_{op}} (1 - \varepsilon)$$

□