

"BENIGN OVERFITTING IN LINEAR REGRESSION"

- BARTLETT, LONG, LUGOSI, TSIGLER, 2019.

BACKGROUND:

Parametric family of fcts:  $F = \{f(\cdot, \theta) : \theta \in \mathbb{R}^p\}$

Data:  $\{(y_i, \alpha_i)\}_{i \in [m]}$        $\alpha_i \in \mathbb{R}^d$        $y_i \in \mathbb{R}$   
 $\alpha_i \sim_{iid} P$ ,       $y_i = f(\alpha_i, \theta_*) + \varepsilon_i$  ← noise e.g.  $\varepsilon_i \sim_{iid} N(0, \sigma_\varepsilon^2)$

Goal: fit the data with  $F$  i.e. estimate  $\hat{\theta}$  such that  $f(\cdot, \hat{\theta})$  is a 'good prediction' on new data points

e.g.  $\mathbb{E}_{y_{new}, \alpha_{new}} [(y_{new} - f(\alpha_{new}, \hat{\theta}))^2]$  small. ←

$\alpha_{new} \sim P$        $y_{new} = f(\alpha_{new}, \theta_*) + \varepsilon_{new}$

Classical approach:

$$\hat{\theta}(\lambda) = \underset{\theta \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ \underbrace{\sum_{i=1}^m (y_i - f(\alpha_i, \theta))^2}_{\text{fit error}} + \lambda \|\theta\|_2^2 \right\}$$

Test error:  $R(\theta_*, \lambda) = \mathbb{E}_\alpha [(f(\alpha, \theta_*) - f(\alpha, \hat{\theta}(\lambda)))^2]$  ←

prediction on new data point:      noise  $\varepsilon \sim N(0, \sigma_\varepsilon^2)$        $y_{new} = f(\alpha_{new}, \theta_*) + \varepsilon_{new}$

$$\mathbb{E}_{\alpha_{new}, y_{new}} [(y_{new} - f(\alpha_{new}, \hat{\theta}))^2] = \underbrace{R(\theta_*, \lambda)}_{\text{test error}} + \sigma_\varepsilon^2 \quad \checkmark$$

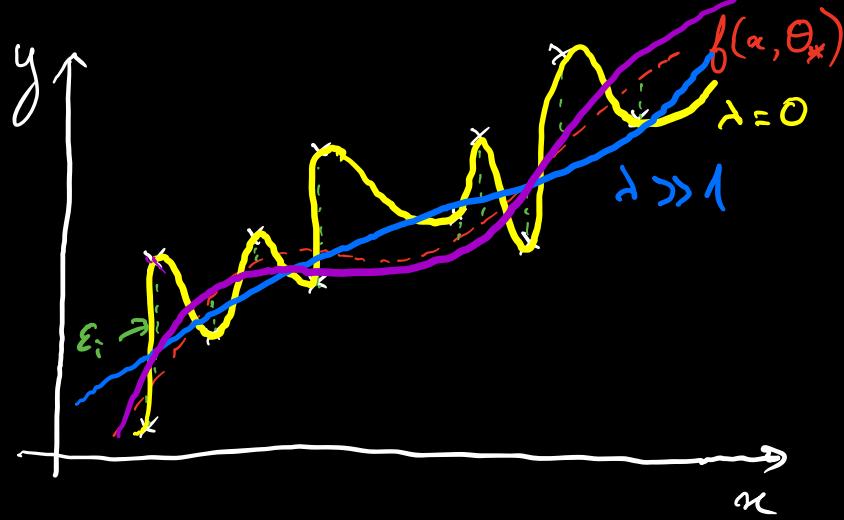
Bias - Variance decomposition:  $y_i = f(\alpha, \theta_*) + \varepsilon_i$

more  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_m) \rightarrow \hat{\theta}(\lambda)$  depends on  $\varepsilon$

$$\begin{aligned}\mathbb{E}_{\varepsilon}[R(\theta, \lambda)] &= \mathbb{E}_{\alpha, \varepsilon} \left[ (f(\alpha, \theta_*) - f(\alpha, \hat{\theta}(\lambda)))^2 \right] \\ &= \mathbb{E}_{\alpha} \left[ (f(\alpha, \theta_*) - \mathbb{E}_{\varepsilon} [f(\alpha, \hat{\theta}(\lambda))])^2 \right] + \mathbb{E}_{\alpha, \varepsilon} \left[ (f(\alpha, \hat{\theta}(\lambda)) - \mathbb{E}_{\varepsilon} [f(\alpha, \hat{\theta}(\lambda))])^2 \right] \\ &= \text{BIAS } (\lambda) + \text{VARIANCE } (\lambda) \\ &= \mathbb{E}_{\alpha} \left[ \left\{ \text{bias} (\mathbb{E}_{\varepsilon} [f(\alpha, \hat{\theta})]) \right\}^2 \right] &= \mathbb{E}_{\alpha} \left[ \text{Var}_{\varepsilon} (f(\alpha, \hat{\theta})) \right]\end{aligned}$$

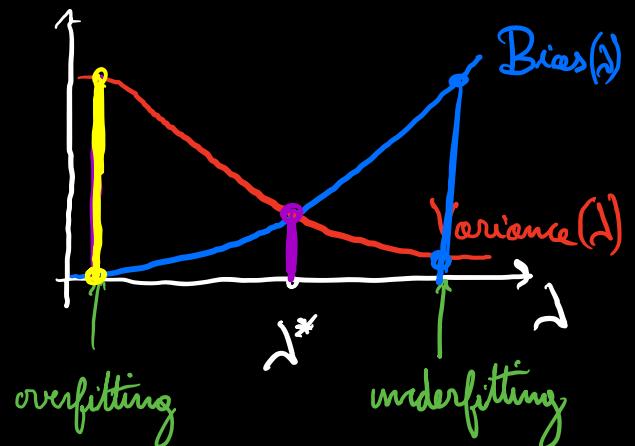
CLASSICAL PICTURE:

$\lambda = 0$        $\lambda$  very large



$$\hat{\theta}(\lambda) = \underset{\theta \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ ER(\theta) + \lambda \|\theta\|_2^2 \right\}$$

$$F \left\{ \|\theta\| \leq C(\lambda) \right\}$$



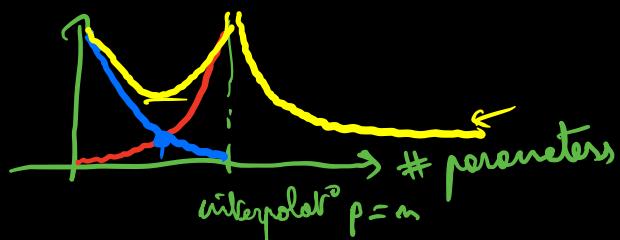
Trade-off between BIAS and VARIANCE

→ need to carefully choose regularization parameter  $\lambda$  in order to control the complexity of the fitted model  $f(\alpha, \hat{\theta}(\lambda))$

⇒ too complex: small bias, large variance (overfitting)

⇒ too simple: large bias, small variance (underfitting)

## MODERN APPROACH:



Train until interpolation:  $f(x_i, \hat{\theta}) = y_i = f(x_i, \theta_*) + \epsilon_i$

→ IMPLICIT REGULARIZATION: choose interpolator ↪

but no explicit control on the bias-variance trade-off  
(no control on the complexity of the function)

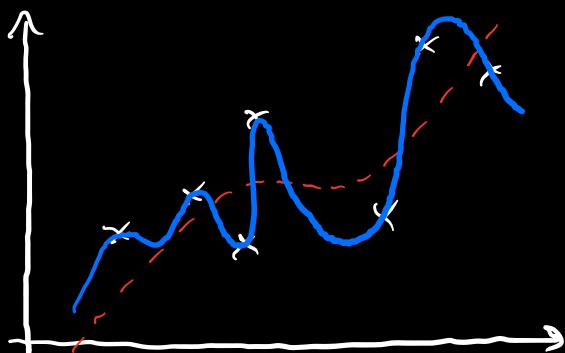
HERE: I will present a scenario where the balance between the bias and the variance is achieved not by carefully tuning an explicit parameter but by a novel phenomena:

### SELF - INDUCED REGULARIZATION

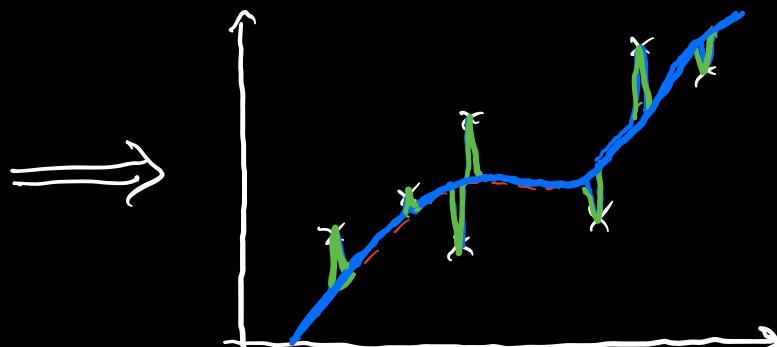
≠ implicit regularization: choose an interpolator  $\hat{\theta}$

given an interpolator  $\hat{\theta}$ : self-induced regularization shows that interpolation does not hurt generalization (test error)

### Example of BENIGN OVERFITTING PHENOMENA



Expected picture



Sometimes "benign"

## Benign overfitting:

$$\text{estimator} : \hat{f}(x, \hat{\theta}) = \underbrace{\hat{f}_0(x, \hat{\theta}_0)}_{\substack{\text{good for prediction} \\ \text{because smooth}}} + \Delta(x)$$

spikes  $\downarrow$   
that are  
useful for interpolation  
but do not harm prediction  
because  $\mathbb{E}_{\alpha}[\Delta(\alpha)^2] < 1$ .

# RIDGE REGRESSION IN LINEAR MODELS:

- Setting:
- Data:  $\{(y_i, \alpha_i)\}_{i \leq m}$   $\alpha_i \in \mathbb{R}^p$ ,  $y_i \in \mathbb{R}$
  - $\mathbb{E}[\alpha_i] = 0$   $\mathbb{E}[\alpha_i \alpha_i^T] = \Sigma$  ↗  
 $\underline{\alpha_i = \sum z_i}$  where  $z_i$  is c-sub-Gaussian
  - $\alpha_i$  are iid,  $y_i = \langle \Theta_*, \alpha_i \rangle + \varepsilon_i$   
 noise  $\varepsilon_i$  independent  $\mathbb{E}[\varepsilon_i] = 0$ ,  $\mathbb{E}[\varepsilon_i^2] = \sigma^2$

Fit a linear model:  $f(x, \theta) = \langle x, \theta \rangle$   $\theta \in \mathbb{R}^p$

Ridge regression:

$$\hat{\theta}_\lambda(x, y) = \underset{\theta \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ \|y - x\theta\|_2^2 + \lambda \|\theta\|_2^2 \right\}$$

$$x = \begin{bmatrix} & \xrightarrow{p} \\ \uparrow & \alpha_1 \\ \vdots & \alpha_m \end{bmatrix} \quad y = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix}$$

$$= \underline{x^T S y} \quad \text{where } \underline{S = (x x^T + \lambda \text{Id})^{-1}}$$

Test error:  $R(\hat{\theta}, x, y, \lambda) = \mathbb{E}_{\alpha} [ (f(x, \theta_*) - f(x, \hat{\theta}_\lambda(x, y)))^2 ]$

$$= \mathbb{E}_{\alpha} [ \langle \alpha, \theta_* - \hat{\theta}_\lambda(x, y) \rangle^2 ]$$

$$= \frac{\|\hat{\theta}_\lambda(x, y) - \theta_*\|_{\Sigma}^2}{\sum} = \mathbb{E} \left[ \langle \hat{\theta}_\lambda(x, y), \frac{\alpha \alpha^T (\theta_* - \theta)}{\sum} \rangle \right]$$

$$\left( \|\theta\|_{\Sigma}^2 := \langle \theta, \Sigma \theta \rangle \right)$$

$$\begin{aligned}
 R(\hat{\theta}, X, \lambda) &= \mathbb{E}_{\varepsilon} \left[ \|\hat{\theta}_*(X, \varepsilon) - \theta_*\|_{\Sigma}^2 \right] \\
 &= \mathbb{E}_{\varepsilon} \left[ \|\underbrace{X^T S X \theta_*}_{\text{green}} + \underbrace{X^T S \varepsilon}_{\text{green}} - \theta_*\|_{\Sigma}^2 \right] \\
 &= B(X, \lambda) + V(X, \lambda) \quad \text{green}
 \end{aligned}$$

where  $B(X, \lambda) = \|(\text{Id} - X^T S X) \theta_*\|_{\Sigma}^2$

$V(X, \lambda) = \sigma^2 \text{Tr}(S X \Sigma X^T S)$

### SELF-INDUCED REGULARIZATION:

WLOG: assume  $\Sigma = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$

$$\lambda_1 \geq \lambda_2 \geq \dots$$

Consider top  $k$  eigenspace (corresponding top  $\lambda_1 \dots \lambda_k$ )

write  $\alpha^T = [\underbrace{\alpha_0^T}_{k \text{ first coordinates}}, \underbrace{\alpha_+^T}_{p-k \text{ last coordinates}}]$

$$X = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} = \frac{1}{m} \begin{bmatrix} \xrightarrow{k} & \xleftarrow{p-k} \\ X_0 & | & X_+ \end{bmatrix}$$

$$\begin{aligned}
 \Sigma &= \frac{1}{n} \left( \begin{array}{c|c} \Sigma_0 & \\ \hline & \Sigma_+ \end{array} \right) & \theta &= (\theta_0, \theta_+) \\
 && \downarrow & \downarrow \\
 && (\theta_1, \dots, \theta_k) & (\theta_{k+1}, \dots, \theta_p)
 \end{aligned}$$

$$XX^T = \underbrace{X_0 X_0^T}_{\text{imagine}} + \underbrace{X_+ X_+^T}_{\downarrow M} \quad \left( \frac{\max(X_+ X_+^T)}{\min(X_+ X_+^T)} \leq L \right)$$

$$\begin{bmatrix} X_0 & X_+ \end{bmatrix} \begin{bmatrix} X_0^T \\ X_+^T \end{bmatrix} \text{ with } \frac{1}{c} \text{Id} \leq M \leq c \text{Id}$$

Recall  $\hat{\Theta}_{\gamma=0}(X, y) = X^T S y \approx X^T (X_0 X_0^T + \gamma \text{Id})^{-1} y + \text{Id}$

$k^{th}$  coordinate:  $\hat{\theta}_0 = \underbrace{X_0^T (X_0 X_0^T + (\gamma) \text{Id})^{-1} y}_{\text{the same as}}$

$\rightarrow$  the same as  $\langle \alpha_0, \Theta_0 \rangle \quad (\lambda + \gamma)$

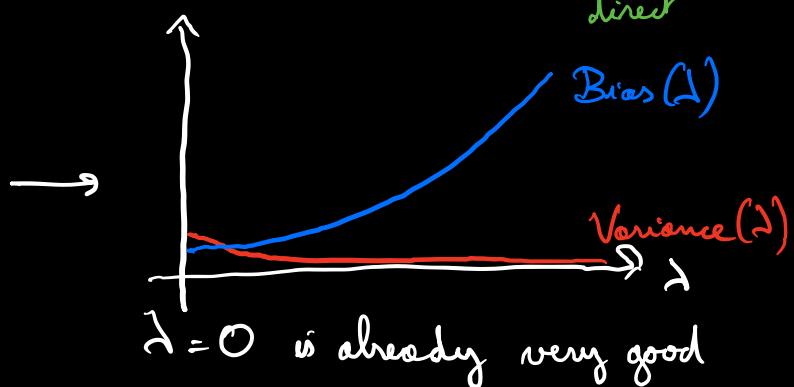
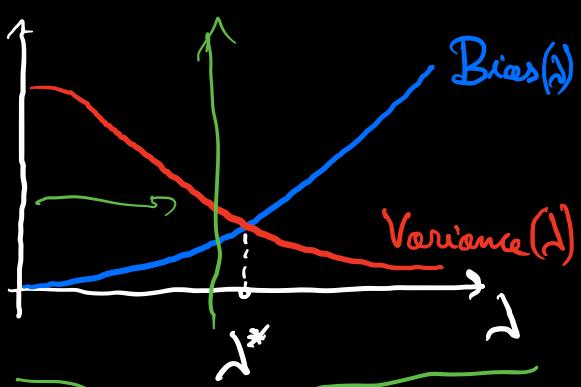
$$\hat{\Theta}_0 = \underset{\Theta_0 \in \mathbb{R}^k}{\operatorname{argmin}} \left\{ \|y - X_0 \Theta_0\|_2^2 + \gamma \|\Theta_0\|_2^2 \right\}$$

*self-induced regularization*

Heuristic:  $X_0$  "signal part" of the features  
 $X_+$  "noise part" of the features

$\rightarrow$  noise part act as an effective ridge regularization

Effective regularization:  $\lambda \rightarrow \lambda_{\text{eff}} = \lambda + \gamma = \sum_{i>n} \lambda_i$   
*low variance direction*



## Main result:

Recall:  $S = (\lambda \text{Id} + X X^T)^{-1} = (\lambda \text{Id} + X_+ X_+^T + X_0 X_0^T)^{-1}$

$$S_+ = \underbrace{(\lambda \text{Id} + X_+ X_+^T)^{-1}}_{\gamma M} \quad S = (S_+^{-1} + X_0 X_0^T)^{-1}$$

Theorem [Bartlett et al., '19] Assume there exists  $k$ :

Assume  $\frac{\lambda_{\max}(S_+)}{\lambda_{\min}(S_+)} \leq L$  with high probability (whp)

Then whp:

$$\text{Bias } (\lambda) \lesssim L^4 \left[ \|\theta_0^*\|_{\Sigma_0^{-1}}^2 \left( \frac{\lambda + \sum_{i>k} \lambda_i}{m} \right)^2 + \|\theta_+^*\|_{\Sigma_+}^2 \right]$$

$$\text{Variance } (\lambda) \lesssim \sigma_\epsilon^2 L^2 \left[ \frac{k}{m} + \frac{m \sum_{i>k} \lambda_i^2}{(\lambda + \sum_{i>k} \lambda_i)^2} \right]$$

Risks: \* Matching lower bounds up to multiplicative constants.

\* In particular, interpolators  $\lambda=0$  are near optimal.

\* Does not fit  $\theta_+^*$  at all:

prediction  $f(x, \hat{\theta}) = \underbrace{\langle \hat{\theta}_0, x_0 \rangle}_{\text{prediction component}} + \underbrace{\langle \hat{\theta}_+, x_+ \rangle}_{\text{interpolation component}}$

$\Rightarrow$  fit well  $\theta_0^*$

\* Small variance: need to choose  $k$  so that

$$(1) \quad k \ll m$$

$$\cancel{(2)} \quad \frac{\left(\sum_{i>m} \lambda_i\right)^2}{\sum_{i>m} \lambda_i^2} \gg m$$

E.g.  $d_{k+1} := \# \text{ eigenvalues of } \sum_{i>m} \left[ \frac{\lambda_{k+1}}{\epsilon}, \lambda_{k+1} \right]$

and assume eigenvalues  $< \frac{\lambda_{k+1}}{\epsilon}$  are negligible.

$$\Rightarrow (2) \text{ requires } \frac{\left(d_{k+1} \cdot \lambda_{k+1}\right)^2}{d_{k+1} \cdot \lambda_{k+1}^2} \gg m \Rightarrow d_{k+1} \gg m$$

and we have

$$\text{Variance} \approx \frac{k}{m} + \frac{m}{d_{k+1}}$$

variance of fitting the rest  
effect of overfitting  
↓  
negligible

Variance of fitting  $\Theta_0^*$   
parametric rate of fitting  $\Theta_0^*$

TO GET BENIGN OVERFITTING HERE:  $\sum p \gg m$

1) Overparametrization: number of directions in parameter space that are unimportant for prediction  $\gg$  sample size

2) Smallest eigenvalues of  $\Sigma$  must decay slowly.

$$d_{k+1} \gg m$$

Self-induced regularization: (if  $\alpha$  satisfies some small ball properties)

Lemma:  $\exists c > 0$  s.t. with prob at least  $1 - 2 \exp(-\frac{m}{c})$

$$\frac{\lambda_{\min}(X, X^T)}{\lambda_{\min}(X_+, X_+^T)} \leq c \left( \frac{\sum_{i \geq k} \lambda_i + m \lambda_{k+1}}{\sum_{i \geq k} \lambda_i - c m \lambda_{k+1}} \right) \leq L$$

Properties of  $\Sigma$  are crucial in this setting.

$$m \leq \frac{1}{2c} \left( \frac{\sum_{i \geq k} \lambda_i}{\lambda_{k+1}} \right)$$

$$m \ll \frac{\lambda_{k+1} \lambda_{k+1}}{\lambda_{k+1}}$$

Andree's lecture: Kernel Ridge Regression

$$K = (h(\langle \alpha_i, \alpha_j \rangle))_{i,j \in [m]}$$

$$\text{fit a fct: } \hat{f}_\lambda(\alpha; \hat{\alpha}) = \sum_{i=1}^m \hat{\alpha}_i (\gamma) h(\langle \alpha_i, \alpha \rangle)$$

$$\text{where } \hat{\alpha} = (K + \lambda \text{Id})^{-1} y$$

Here: self induced regularization comes from non polynomial part of the kernel fct:

$$h(x) = \sum_{k=0}^{\infty} \frac{h^{(k)}(0)}{k!} x^k$$

$$K = \sum_{k=0}^{\infty} \frac{h^{(k)}(0)}{k!} Q_k \quad \text{where } Q_k = (\langle \alpha_i, \alpha_j \rangle^k)_{i,j \in [n]}$$

Then if  $d^l \ll m \ll d^{l+1}$ ,

$$\left\| \sum_{k=l+1}^{\infty} \frac{h^{(k)}(0)}{k!} Q_k - \gamma \text{Id} \right\|_{\text{op}} = O_d(1).$$

$$K \approx K_{\leq l} + \underbrace{\gamma \text{Id}}_{\text{self-induced regularization}}$$

$\Rightarrow$  leads also to benign overfitting.

[Bartlett et al] result: applies also to  $p = \infty$

$\rightarrow \alpha$ : could be the vector of an orthonormal basis of  $h$ .

$\rightarrow$  however not subgaussian: requires very different techniques

Nonetheless:

Kernel: eigenvalues

$$\lambda_k \approx d^{-k} \text{ with degeneracy } O_d(d^k)$$

$\hookrightarrow$  associated to polynomials of degree  $k$ .

$\Rightarrow$  to fit all degree  $l$  polynomials:

$$\frac{k}{m} + \frac{m}{d^{l+1}}$$

$$``d_k" = \sum_{s=0}^l d^s \approx d^l$$

$$``d_{l+1}" \approx d^{l+1}$$

$$\text{Variance} \approx \frac{d^l}{m} + \frac{m}{d^{l+1}}$$

$$+ \text{Bias} \approx \| P_{\geq l} f_* \|_2^2$$

$$K(\alpha_i, \alpha_j) = h(\langle \alpha_i, \alpha_j \rangle)$$

$$K_m = (K(\alpha_i, \alpha_j))_{i,j \in [m]}$$

$$h(u) = \sum_{k=0}^{\infty} \frac{h^{(k)}(0)}{k!} u^k$$

$$K_m = \sum_{k=0}^{\infty} \frac{h^{(k)}(0)}{k!} Q_k \quad Q_k = (\langle \alpha_i, \alpha_j \rangle^k)_{i,j \in [m]}$$

$$n \ll d^{l+1}$$

$$\left\| \sum_{k=l+1}^{\infty} \frac{h^{(k)}(0)}{k!} Q_k - \gamma \text{Id} \right\|_{\text{op}} \ll 1$$

$$\hat{\Theta} = (K + \gamma \text{Id})^{-1} y$$

$$\approx (\underline{K_{\leq l}} + (\underline{\Delta} + \gamma) \text{Id})^{-1} y$$

→ able to fit  $P_{\leq l} f_*$  and nothing else

$$f_*(\alpha) = \langle \alpha, \Theta_* \rangle \rightarrow \rho = \infty$$

$$h(\langle \alpha, y \rangle) = \sum_{i=1}^{\infty} \gamma_i \downarrow_i \Phi_i(\alpha) \Phi_i(y)$$

$\downarrow$

$\{\Phi_i\}$  orthonormal basis  
 $L^2(\mathbb{P})$

$$f_*(\alpha) = \sum_{i=1}^{\infty} \sqrt{\gamma_i} \Phi_i(\alpha) \Theta_i^*$$

$$= \langle \alpha, \Theta_i^* \rangle$$

$$\underbrace{\alpha = (\sqrt{\gamma_i} \Phi_i(\alpha))_{i \geq n}}$$

$$Q_k = (\langle \alpha_i, \alpha_j \rangle^k)_{i,j \in [n]} \quad \alpha_i \sim \text{Unif}(\mathbb{S}^{d-1}(1))$$

$$Q_k = \text{Id} + \Delta \rightarrow \text{off diagonal elements of } Q_k$$

$$\langle \alpha_i, \alpha_j \rangle \leq \frac{1}{\sqrt{d}} \text{ w.h.p.}$$

$$\max_{i \neq j} \langle \alpha_i, \alpha_j \rangle \leq \frac{1}{d^{\frac{1}{2}-\varepsilon}}$$

$$\|\Delta\|_F^2 = \sum_{i \neq j} \langle \alpha_i, \alpha_j \rangle^{2k} \leq n^2 d^{-k + \frac{\varepsilon}{2} - 2k} \rightarrow 0$$

$$n \leq d^{\frac{k}{2}} \leftarrow \text{Frobenius}$$

$$n \leq d^{\frac{k}{2} - \delta} \leftarrow \text{complicated}$$

$$d^l \ll n \ll d^{l+1}$$

$\rightarrow$  fit exactly  $P \leq l$  for  $f$

$\rightarrow$  even  $\lambda = 0$

$\lambda \rightarrow \lambda_i > 0$  for any eigenvalue of polynomials  $\leq l$

$\lambda = 0 \rightarrow h$  not a polynomial  $d^0 \leq l$