

# Random Function Classes

①

## Kernel (recall)

$$k(x, x') = \langle \ell(x), \ell(x') \rangle$$

$$= \sum_i \lambda_i \ell_i(x) \ell_i(x')$$

where  $\langle \ell_i, \ell_j \rangle = \delta_{ij}$  (in  $L_2$ )

## Corollary:

Sampling yields a good rate:

draw  $i$  with  $p(i) = \frac{\lambda_i}{\sum_j \lambda_j}$  has

variance  $\frac{1}{m} \left[ \sum_i \lambda_i (\ell_i(x) \ell_i(x'))^2 - k(x, x')^2 \right]$

## Useful Lemmas:

(Neal, 1994)

given  $\ell_2$  with  $z \sim p(z)$

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

yields function  $f(x) := \frac{1}{m} \sum_i \alpha_{z_i} \ell_{z_i}(x)$

such that

$$\mathbb{E}[f(x)] = \mathbb{E}_z \left[ \sum_i \mathbb{E}[\alpha_{z_i}] \ell_{z_i}(x) \right] = 0$$

$$\mathbb{E} \left[ \frac{f(x) f(x')}{f(x) f(x')} \right] = \mathbb{E}_z \left[ \frac{1}{m} \sum_{i=1}^m \mathbb{E}[\alpha_{z_i}^2] \ell_{z_i}^2(x) \right]$$

$$= \mathbb{E}_{z \sim p(z)} [\ell_2(x) \ell_2(x')] =: k(x, x')$$

$\Rightarrow$  This is a Gaussian Process  $\mathcal{G}$

If  $|\ell_2(x)| \leq R$

then  $P_2 \left[ \left| k(x, x') - \frac{1}{m} \sum_{i=1}^m \ell_{z_i}(x) \ell_{z_i}(x') \right| > \epsilon \right]$

$$\leq 2 e^{-m \epsilon^2 / R^4}$$

by Application of Hoeffding  $\Rightarrow \epsilon = O\left(\frac{\sqrt{\log S}}{m}\right)$

Much Stronger result via Matrix Hoeffding (we will prove this later in class)

Then: given matrices  $Y_k$  with  $\mathbb{E}[Y_k] = 0$

and  $Y_k^2 \preceq A_k \quad \forall k$

then

$$P \left\{ \lambda_{\max} \left( \sum_k Y_k \right) \geq \epsilon \right\} \leq d e^{-\frac{\epsilon^2}{26L}}$$

where  $L = \frac{1}{2} \left\| \sum_k A_k + \mathbb{E}[Y_k^2] \right\|$

Application

$$Y_k := \frac{1}{m} [\ell_2(\bar{x}) \ell_2(x) - k]$$

$$\mathbb{E}[Y_k] = 0$$

$$Y_k^2 = \frac{1}{m^2} \left[ \underbrace{\ell_2 \ell_2^+ \ell_2 \ell_2^+}_{4 \ell_2 \ell_2^2} - \ell_2 \ell_2^+ k - k \ell_2 \ell_2^+ \right]$$

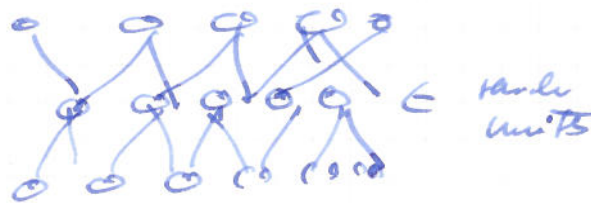
assume that  $\|\ell_2\|^2 \leq m R^2$

then we pick  $A_k = m R^2 \cdot \mathbb{1}$

# Random Function Classes

This has been (re)invented many times:

- a) Neural Networks with random activation potential  
~ 1980s



- b) Empirical Kernel Map (Toschi, Schölkopf, ~ 1990)  
 $x \rightarrow (k(x_1, x), \dots, k(x_n, x))$

$$\Rightarrow \bar{k}(x, x') = \int_{x \sim p(x)} [k(x, x') k(x', x)]$$

- c) Random Kitchen Sinks (Rahimi & Recht, ~ 2008)

for Gaussian / RBF kernel

$$k(x - x') = \int d\mathcal{G}(\omega) \frac{e^{i\langle \omega, x - x' \rangle}}{e^{i\langle \omega, x \rangle} e^{-i\langle \omega, x' \rangle}}$$

$$f_\omega(x) := (\sin(\omega_1 x), \cos(\omega_1 x))$$

renormalize measure  $d\mathcal{G}(\omega)$  to get

$$\omega_j \sim \frac{g(\omega)}{\|g(\omega)\|_1}$$

$$\text{and } \hat{k}(x, x') = \frac{1}{n} \sum_{j=1}^n e^{i\langle \omega_j, x \rangle} e^{-i\langle \omega_j, x' \rangle}$$

Recipe: for Gaussian RBF draw  $\Omega \sim \mathcal{N}(0, 1)$

$\frac{1}{\sqrt{2}} \exp(i(\Omega \cdot x))$  are basis functions  
Otherwise simply rescale  $\Omega$ .

- d) inner product (and  $k(x, x') = k(\langle x, x' \rangle)$ )

(2)

Invariant  $k(x, x') = k(Ux, Ux')$

This means that in its representation on  $\ell_2$  we need to have  $[g(u) f(x)]^T [g(u) f(x')] = k(\langle x, x' \rangle)$

Hence  $f$  decomposes into irreducible orthogonal representations  $k(x, x') = \sum_j \lambda_j f_j(x) f_j(x')$

Separating radial / angular parts implies an expansion in terms of spherical harmonics, hence we have  $k(x, x') = \sum_j \lambda_j L_j(\frac{\langle x, x' \rangle}{\|x\| \|x'\|})$  where  $L_j$  are adjoint Legendre polynomials (chosen for  $\mathbb{R}^d$ )

use the fact that

$$L_j(\langle x, x' \rangle) = \int_{S^{d-1}} d\Omega L_j(\langle x, z \rangle) L_j(\langle x', z \rangle)$$

on the unit sphere (i.e.  $\|x\| = \|x'\| = 1$ )

Recipe:  $\ominus$  draw  $\Omega \sim \mathcal{N}(0, 1)$

$\ominus$  rescale rows to unit length

$\ominus$  compute  $L_j(\langle x, z \rangle)$

for  $j \sim p(j)$  according to choice of kernel



# Random Function Classes

e) Fastfood (Sarlos, Le, Smola '13)  
 key bottleneck is computing and storing of  $\Omega$   
 but this is a Gaussian random matrix...

Idea: approximate  $\Omega$  by  $SHT\pi GB$   
 where  $S \in$  Scaling matrix (diagonal)  
 $H \in$  Hadamard  
 $\pi \in$  Permutation  
 $B \in$  Binary digital (diagonal)  
 $G \in$  Gauss (diagonal)

Runtime Complexity: diagonal of  $\pi$  are  $O(n)$

$$H_{2n} = \begin{bmatrix} H_n & H_n \\ H_n & -H_n \end{bmatrix} \text{ can be multiplied in } O(n \log n) \text{ time.}$$

Proof: assume this is possible for  $H_n$   
 $\Rightarrow H_{2n} \begin{pmatrix} v_{2n} \\ v_{2n} \end{pmatrix} = \begin{pmatrix} H_n (v_n + v_n) \\ H_n (v_n - v_n) \end{pmatrix}$

$$\text{this is } O\left(\frac{2 \cdot (n \log n) + 2n}{2n \log 2n}\right)$$

(do not use Matlab's FFT  $\rightarrow$  it's SLOW)

$\Rightarrow O(n \log d)$  CPU,  $O(n)$  RAM



$\rightarrow$  well-definedness:

$$\left[ (HT\pi GB)^T (HT\pi GB) \right]_{ii} = n \cdot \left[ B H^T G^2 H B \right]_{ii} = \underline{n \operatorname{tr} G^2}$$

(all random vectors lie on a sphere)  $\square$

$\rightarrow$  for a given row of  $HT\pi GB$ , we have independent Gaussians:

- binary doesn't change things (Central Gaussian)
- ditto Hadamard  $\checkmark$

$\rightarrow$  rats are a lot of work (maybe cover the device a la Chazelle later).

f) Random Set indicators (Davies & Ghahramani '14)

- $\rightarrow$  generate random set indicators (e.g. trees)
  - $\rightarrow$  compute kern  $k(x, x') = \mathbb{E}[\delta(S(x), S(x'))]$
  - $\rightarrow$  example: random decision trees / forest
- all this is useful for

kernel / function classes (no point in using dual space unless  $n \ll d$ ).