

Lecture 4: infinite width over \mathbb{R}^d ; Fourier transforms & Bernoulli terms.

Announcements:

* Picnic OH (more social than technical)

* HW1 out; have 20 days.
Later HW shorter.

Recap.

* "Universal approximation"; worst-case, non-constructive.

* Infinite-width: univariate case, by FTC

$$\forall x \in [0, 1], \quad g(x) - g(0) = \int_0^x \mathbb{1}_{[x \leq b]} g'(b) db$$

* We'll say more about sampling today.

Also today: Fourier & infinite-width over \mathbb{R}^d .

* Why infinite-width

* Used in many optimization analyses.

* Removes/simplifies many error terms.

* In early phase of training, weights

near random initialization; behaves like

infinite-width due to concentration.

Sampling

Given $\int_0^1 \sigma(x-b) g(b) db = g(x)$

1. Sample $(b_1, \dots, b_m) \sim \frac{|g(b)|}{\int_0^1 |g(b)| db}$.
2. For each j , set $a_j := \text{sgn}(g(b_j)) \cdot \int_0^1 |g(b)| db \cdot \frac{1}{m}$
3. Output $f(x) := \sum_j a_j \sigma(x-b_j)$

Unbiased $\mathbb{E}_{b_1, \dots, b_m} f(x) = g(x)$

Error estimate:

$$\int_0^1 (f(x) - g(x))^2 dx \leq \frac{1}{m} \left(\int_0^1 |g(b)| db \right)^2 \max_{b \in [0,1]} \int_0^1 \sigma(x-b)^2 dx$$

\uparrow \uparrow \uparrow
 \downarrow or \uparrow number nodes mass of weights penalty for activation

Remark (\mathbb{R}^d)

* in notes; $\frac{1}{m}$ still there (not $\frac{1}{m^{1/d}}$),

'mass' harder to estimate, might $\Omega(d)$

* More general version in notes.

Fourier & infinite width (Barron '93)

Goals (for multivariate infinite-width representation)

$$* f(x) = \int \sigma(w^T x - b) g(b, w) db dw$$

* Mass $\int |g(b, w)| db dw$ to not be worst case;

e.g., $\int |g(b, w)| db dw = \text{poly}(d)$ for some reasonable cases.

Recall: Fourier transform

$$\hat{f}(w) = \int \underbrace{\exp(-2\pi i w^T x)}_{\text{activation}} f(x) dx.$$

Fourier inversion: $f, \hat{f} \in L_1$ ^{step down}

$$f(x) = \int \exp(2\pi i w^T x) \hat{f}(w) dw,$$

as in finite-width network!

Problem: want typical activation, not complex sinusoids.

Remark:

* This infinite-width form is not in (Barron '93), but first few steps of derivation are.

to start deriving threshold form, since $\text{Re}[z] = f$,

$$f(x) = \text{Re}[f(x)] = \text{Re}\left[\int \exp(2\pi i w^T x) \hat{f}(w) dw\right]$$

$$= \int \text{Re}\left[\exp(2\pi i w^T x) \hat{f}(w)\right] dw$$

polar form $\hat{f}(w) = \underbrace{|\hat{f}(w)|}_{\text{Magnitude}} \underbrace{\exp(2\pi i \theta(w))}_{\text{orientation}}$

$$= \int \text{Re}\left[\exp(2\pi i (w^T x + \theta(w))) \underbrace{|\hat{f}(w)|}_{\text{real}}\right] dw$$

$$= \int \underbrace{\text{Re}\left[\exp(2\pi i (w^T x + \theta(w)))\right]}_{\text{real-valued activation!}} \cdot |\hat{f}(w)| dw$$

$$e^{iz} = \cos(z) + i \sin(z)$$

$$\cos(2\pi (w^T x + \theta(w)));$$

note

$$\cos(2\pi (w^T x + \theta(w))) - \cos(2\pi \theta(w))$$

$$= - \int_0^{w^T x} 2\pi \sin(2\pi (b + \theta(w))) db$$

$$= -2\pi \int_0^{\|w\|} \mathbb{1}[w^T x \geq b] \sin(2\pi (b + \theta(w))) db$$

$$- \int_{-\|w\|}^0 \mathbb{1}[-w^T x \geq -b] \sin(2\pi (b + \theta(w))) db$$

$$* f(x) - f(0) = -2\pi \left[\int_0^{\|w\|} \mathbb{1}[w^T x \geq b] \sin(2\pi (b + \theta(w))) |\hat{f}(w)| db dw - \int_{-\|w\|}^0 \mathbb{1}[-w^T x \geq -b] \sin(2\pi (b + \theta(w))) |\hat{f}(w)| db dw \right]$$

$$\iint \mathbb{1}[w^T x \geq b] g(b, w) db dw$$

mass:

$$2 \cdot 2\pi \iint_0^{\|w\|} |\sin(2\pi (b + \theta(w)))| \cdot |\hat{f}(w)| \cdot db dw$$

$$\leq 2 \cdot 2\pi \int_0^{\|w\|} db |\hat{f}(w)| dw$$

$$= 2 \cdot 2\pi \cdot \int \|w\| \cdot |\hat{f}(w)| dw$$

$$= 2 \cdot \int \|\nabla \hat{f}(w)\| dw.$$

fact: $\|\nabla \hat{f}(w)\| \geq 2\pi \|w\| \cdot |\hat{f}(w)|$

Definition:

"Barron norm"

"Barron class"

$$\{ f : \int \|\nabla \hat{f}(w)\| dw < C \}$$

Theorem: If $f, \hat{f} \in L_1$,

$$f(x) - f(0) = -2\pi \left[\int_0^{\|w\|} \mathbb{1}[w^T x \geq b] \sin(2\pi (b + \theta(w))) |\hat{f}(w)| db dw - \int_{-\|w\|}^0 \mathbb{1}[-w^T x \geq -b] \sin(2\pi (b + \theta(w))) |\hat{f}(w)| db dw \right]$$

$$\text{Weight mass} \leq 2 \int \|\nabla \hat{f}(w)\| dw.$$

Remarks:

* Some Fourier mass estimates.

* $\Omega(2^d)$ for many radial functions.

* "almost flat" Gaussians have $\int \|\nabla \hat{f}(w)\| = O(\sqrt{d})$

* Bound NTK representation size using

Barron norms.