

# Lec 7: "kernel" part of NTK; start of deep/shallow separation.

## Announcements

- \* HW1 due 9/22 (no late submissions!)
- \* NTK notes up; section 5 tonight.
- \* Loomis HW Thursday only; short OH today.

Question.

apx		} classed view
opt		
gen	xx	

function class increase

DL adapts to complexity of the problem

some classical theory blows up as  $m \rightarrow \infty$

## Recap

$$f(x; W) := \frac{1}{\sqrt{m}} \sum_{i=1}^m a_i \sigma(w_i^T x)$$

$$f_0(x; W) := f(x; W_0) + \langle \nabla f(x; W_0), W - W_0 \rangle$$

$$|f(x; W) - f_0(x; W)| \leq \begin{cases} O\left(\frac{\|W - W_0\|_F^2}{\sqrt{m}}\right) & \sigma \text{ smooth} \\ O\left(\frac{\|W - W_0\|_F^{4/3}}{m^{1/2}}\right) & \sigma \text{ ReLU} \\ = O\left(\left(\frac{\|W - W_0\|_F}{\sqrt{m}}\right)^{4/3}\right) \end{cases}$$

Today: "kernel story".

$$k_m(x, x') := \langle \nabla f(x; W_0), \nabla f(x'; W_0) \rangle$$

$$= \left\langle \begin{bmatrix} a_1 \mathbb{1}[w_{0,1}^T x \geq 0] x / \sqrt{m} \\ \vdots \\ a_m \mathbb{1}[w_{0,m}^T x \geq 0] x / \sqrt{m} \end{bmatrix}, \begin{bmatrix} x' \\ \vdots \\ x' \end{bmatrix} \right\rangle$$

$$= \left(\frac{1}{\sqrt{m}}\right)^2 \sum_{i=1}^m (x^T x') \mathbb{1}[w_{0,i}^T x \geq 0] \mathbb{1}[w_{0,i}^T x' \geq 0]$$

$$= (x^T x') \left[ \frac{1}{m} \sum_{i=1}^m \mathbb{1}[w_{0,i}^T x \geq 0] \mathbb{1}[w_{0,i}^T x' \geq 0] \right]$$

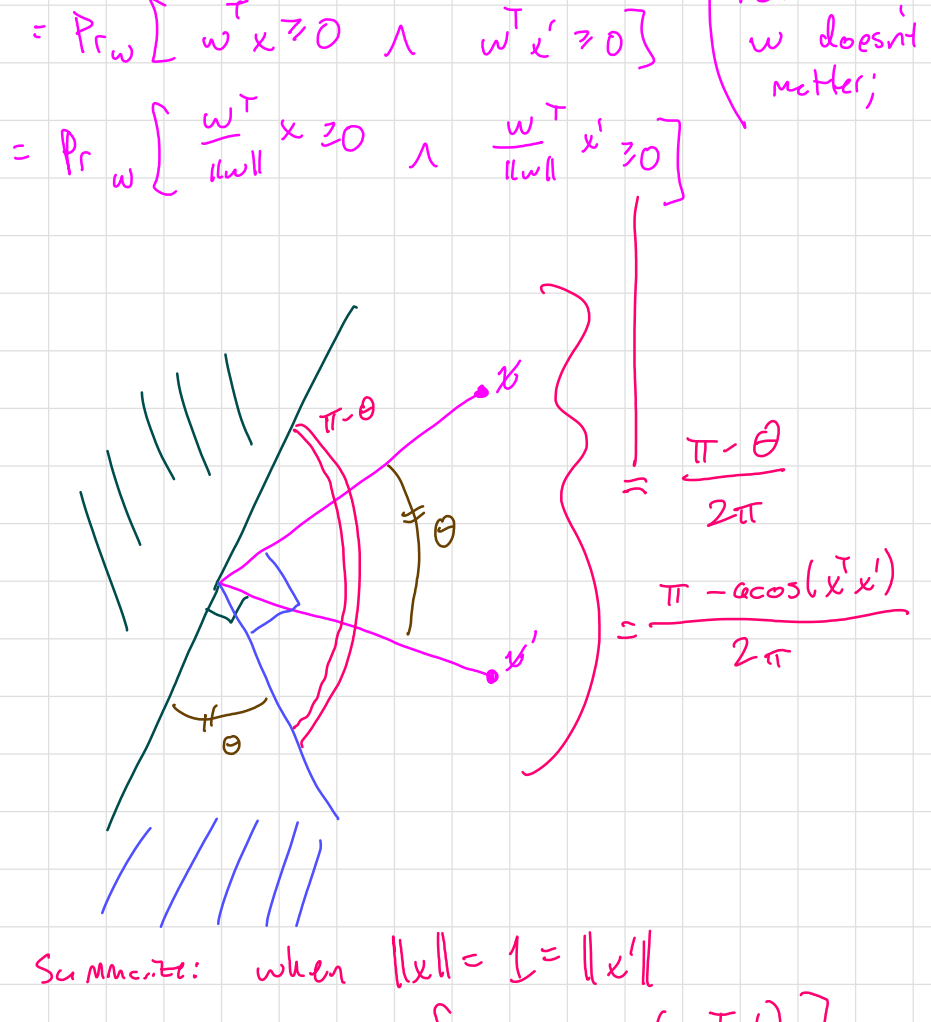
summands are iid.

SLN  $m \rightarrow \infty$   $\rightarrow (x^T x') \mathbb{E}_w \mathbb{1}[w^T x \geq 0] \mathbb{1}[w^T x' \geq 0]$

$k(x, x')$

General activation

 $(x^T x') \mathbb{E}_w \sigma'(w^T x) \sigma'(w^T x')$



Summarize: when  $\|x\| = 1 = \|x'\|$

$$k_{\text{ReLU}}(x, x') = (x^T x') \left[ \frac{\pi - \arccos(x^T x')}{2\pi} \right].$$

## Remark (why kernel):

- Things people do w/ kernel view:
  - \* Evolution of weights can be modeled as "Gaussian Process" (cf. NTK paper Clement-Jacot-Hongler)
  - \* arcos kernel universal approximator.
  - \* Minimum eigenvalue of Gram matrix can lead to optimization rates.

Remark inner product kernels have many special properties which NTK inherits.

Remark. Multi-layer kernel (not well understood) (power vs shallow not well understood)

$$\langle \nabla f(x; \vec{W}(0)), \nabla f(x'; \vec{W}(0)) \rangle$$

$$= \sum_{j=1}^L \langle \nabla_j f(x; \vec{W}(0)), \nabla_j f(x'; \vec{W}(0)) \rangle.$$

Remark. Taylor expansion at 0:

$$\langle \nabla f(x; 0), \nabla f(x'; 0) \rangle$$

$$= \left\langle \begin{bmatrix} a_1 \frac{1}{\sqrt{m}} x \sigma'(0) \\ \vdots \\ a_m \frac{1}{\sqrt{m}} x \sigma'(0) \end{bmatrix}, \begin{bmatrix} a_1 \frac{1}{\sqrt{m}} x' \sigma'(0) \\ \vdots \\ a_m \frac{1}{\sqrt{m}} x' \sigma'(0) \end{bmatrix} \right\rangle$$

using deterministic fixed  $\sigma'(0)$  "semi-gradient"

$$= x^T x' \frac{1}{m} \sum_j a_j^2 \sigma'(0) \sigma'(0) = x^T x' \sigma'(0)^2.$$

Brooks problem:

$$\sigma(\sigma(x) - \sigma(-x)) - \sigma(-x)$$

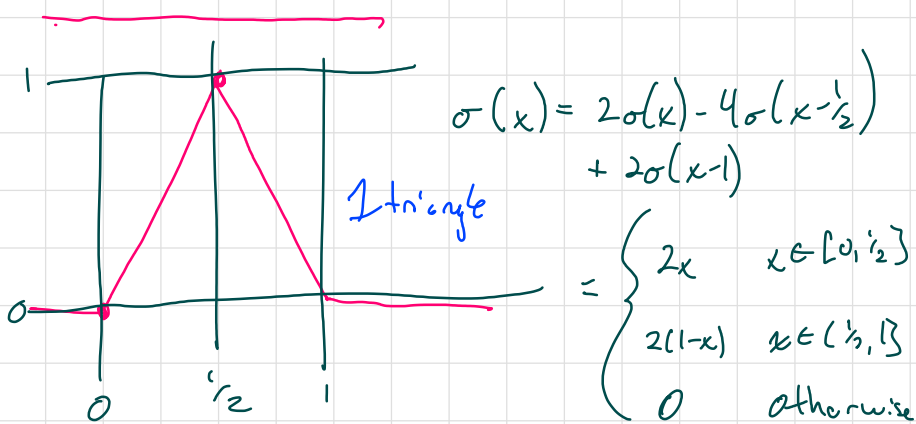
Remark. Can prove arcos kernel gives a universal approximator. (Notes used a bias)

$\Rightarrow$  "closes the loop"

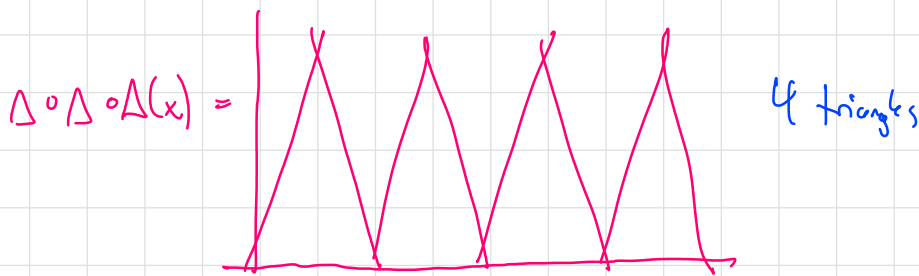
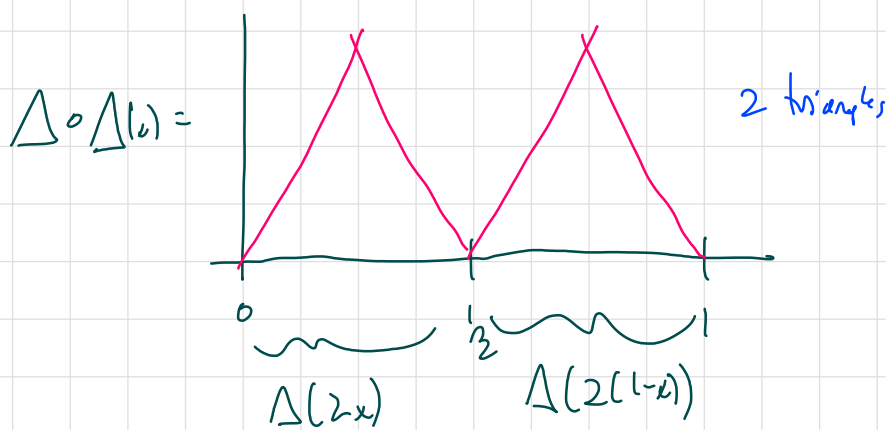
with large  $m$ :  $f_0(\cdot; W) \approx f(\cdot; W)$  continuous  $\mathcal{B}$

# Deep vs shallow (impractical)

- ① " $\Delta$ " quickly increases complexity with depth; only deep vs shallow.
- ②  $\exists$  small deep network which can be approximated by subexponentially-sized shallow networks.
- ③ Can use  $\Delta$  to build: polynomials, Taylor approx...



Credit? (dynamical systems? Bengio et al. "folding"? Turing completeness of LNNs? ...)



$2^{k-1}$  triangles.

$$\Delta^k(x) = (\underbrace{\Delta \circ \dots \circ \Delta}_{k \text{ copies}})(x) = \text{M-M-M}$$

Point: depth allows complexity to increase exponentially fast.

Remark (impractical).

No evidence this relates to power of deep networks trained by gradient descent.

# Office hours

↙  $\mathcal{H}$

$$\left\{ x \mapsto \sum_{j=1}^N \alpha_j k(x, x_j) : \begin{array}{l} N \geq 0, \\ \alpha_j \in \mathbb{R} \\ x_j \in \underline{X} \end{array} \right\}$$

norm:

$$\boxed{\sum_{i,j} \alpha_i \alpha_j k(x_i, x_j)}$$

$$\mathcal{H}_B := \{ f \in \mathcal{H} : \|f\|_{\mathcal{H}} \leq B \}$$

$\forall \varepsilon > 0$  continuous  $g \quad \exists h \in \mathcal{H}$   
s.t.  $\|g - f\| \leq \varepsilon$ .

$$\|f - f_0\| \xrightarrow{m \rightarrow \infty} 0.$$

Given cont  $g \quad \varepsilon > 0$ ,

Pick  $h \in \mathcal{H}, \quad \|h\|_{\mathcal{H}} < \infty$

$$\|f - f_0\| = \frac{\|w - w_0\|_{\mathcal{H}}^2}{\sqrt{m}} \approx \|h\|_{\mathcal{H}}$$

$$\frac{d}{dx} f(rx) = x f'(rx)$$

$$f(z) = z^2 \quad f'(z) = 2z$$

$$\left( \frac{d}{dx} f(rx) = x f'(rx) \right) \Rightarrow 2rx^2$$