



# THEORY OF MACHINE LEARNING

## LECTURE 21

GRADIENT DESCENT FOR NN: NEURAL TANGENT KERNEL

# ANNOUNCEMENTS

- Homework 3 due on Monday April 11

■ Ideas

■ Implicit regularization

■ Stability in optimization

■ Depth vs width (converse)

■ Learning parities - earliest lower bound for deep learning

(possible via other algos; not via deep learning).  
efficient ( $n^3$ ).

■ Last lecture's discussion on feature learning

$$Ax \approx b$$

$$\min \|x\|^2$$

s.t.  $Ax = b$

$$F_3 \text{ vs. } F + \frac{h}{\exp(k)}$$

$$\min \|Ax - b\|^2$$

→ Soln found by grad descent  
will be "simple"

→ This is one reason these  
models don't overfit.

# NEURAL NETWORK TRAINING

- **Question (supervised learning):** given data  $\overset{\substack{\mathbb{R}^d \\ \uparrow}}{(x_1, y_1), (x_2, y_2), \dots}$  from some distribution  $D$ , find  $h$  (with given "architecture") that minimizes the risk

- Really hard theoretically (even if inputs Gaussian and risk zero is achievable)↑

- In practice, solved via gradient descent

[fast implementation (backprop) by Rumelhart, Hinton, Williams]



$$L(w) = \sum_{i=1}^N \left( f(w; x_i) - y_i \right)^2$$

$w^{(0)} \rightarrow w^{(1)} \rightarrow \dots$

- Question for today:

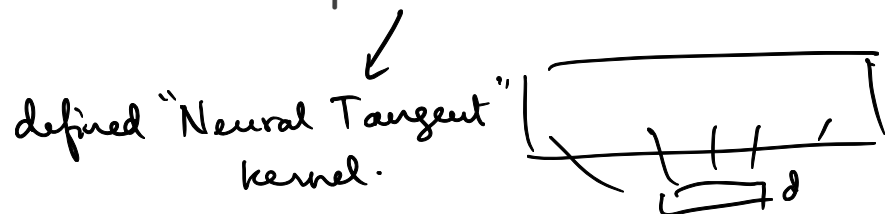
- How can one analyze dynamics of gradient descent?
- Are there cases where we can reason about resulting solution?

# OVER-PARAMETRIZATION

$$\frac{\# \text{parameters}}{|W|} \gg \# \text{inputs } (N).$$

- **Observation:** modern deep nets pretty overparametrized, but they still don't overfit [Why? ... Belkin et al., "double descent"].
- Question out of desperation: Is GD easier to analyze when network is "heavily" overparametrized?  $\rightarrow$  this is one setting where  $\mathcal{L}(\hat{w}) = 0$  [Allen-Zhu, ...].

**Theorem.** [Jacot, Gabriel, Hongler 18] [Arora, et al. 2019] A width  $\sim N^3$  network (any number of layers) trained via GD from random initialization achieves zero training error. Moreover, the final solution is equivalent to solving a "Kernel regression" problem with a specific kernel.



## ASIDE: KERNEL REGRESSION

$(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

- Ubiquitous motivation: function value known at a bunch of points, "interpolate" to rest of space

- One way to think of all of ML!

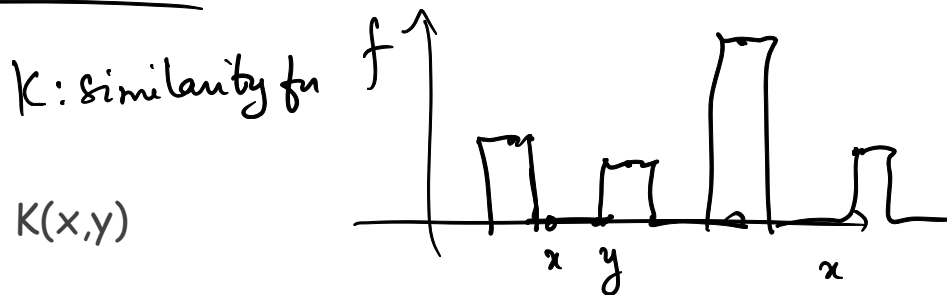
- Suppose  $K$  defines "point similarity"  $K(x, y)$

- Consider interpolation via functions of specific form...

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

Goal: Given kernel (Gaussian), find the best  $\alpha_i$

$\alpha_1 \times \text{similarity with given datapoint } x_1 + \alpha_2 \times \text{similarity with } x_2 + \dots$



$K(x, y)$ : "kernel"

$$K(x, y) = \begin{cases} 1 & \text{if } \|x - y\| < r \\ 0 & \text{otherwise.} \end{cases}$$

$\downarrow$

$$e^{-\frac{\|x - y\|^2}{r^2}}$$

Simplest / most natural : choose  $\alpha_i$  so that  $\hat{f}(x_i) = y_i \quad \forall i$   
 $\hookrightarrow N$  constraints

want to find coeffs  $\alpha_1, \alpha_2, \dots, \alpha_N$ .

$$\hat{f}(x_1) = y_1 \iff \alpha_1 \underbrace{K(x_1, x_1)}_{\substack{\downarrow \\ i \quad n}} + \alpha_2 K(x_1, x_2) + \dots + \alpha_N K(x_1, x_N) = y_1$$

$\hookrightarrow$  linear eqn in the  $\alpha$ 's  
 (since we know the  $x_i \dots$ )

$$\begin{matrix} n \\ i \end{matrix} \left[ \begin{matrix} K \\ \alpha \end{matrix} \right] \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_N \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \rightarrow \text{opt soln: } \alpha = \underline{\underline{K^{-1} y}}$$

$\hookrightarrow$  Try to minimize  $g(\alpha) = \| K\alpha - y \|^2$  using gradient descent (with param =  $\alpha$ )

initialize:  $\alpha^{(0)}$

$$\nabla g^{(\alpha)} = 2 \underline{\underline{K(K\alpha - y)}}$$

$$\underline{\underline{\alpha^{(t+1)} = \alpha^{(t)} - \eta \cdot \nabla g(\alpha^{(t)})}}$$

Qn: What happens <sup>to the  $\hat{f}$</sup>  as you do steps of GD?

$$\hat{f}(\alpha^{(t)}; x) = \sum_{i=1}^N \alpha_i^{(t)} K(x, x_i)$$

$$\hat{f}(\alpha^{(t+1)}; x) - \hat{f}(\alpha^{(t)}; x) = -\eta \langle \nabla \hat{f}(\alpha^{(t)}; x), \nabla g(\alpha^{(t)}) \rangle$$

$$= -\eta \langle \nabla \hat{f}(\alpha^{(t)}; x), \sum_{i=1}^N (\hat{f} - y)_i \nabla \hat{f}(\alpha^{(t)}; x_i) \rangle$$

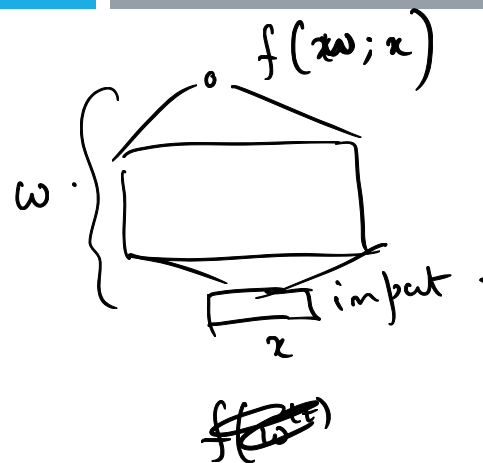
(exercise)

$$= -\eta \sum_{i=1}^N (\hat{f} - y)_i K(x, x_i)$$

---

# NEURAL NET TRAINING REVISITED

$$L(w) = \frac{1}{2} \sum_{i=1}^N (f(w; x_i) - y_i)^2$$



$$\nabla_w L(w) = \sum_{i=1}^N (f(w; x_i) - y_i) \cdot \nabla_w f(w; x_i)$$

$$w^{(t+1)} = w^{(t)} - \eta \cdot \nabla_w L(w^{(t)}) \quad (\text{gradient descent})$$

Qn: how does  $f(w^{(t)}, x)$  "evolve" as time progresses?

$$f(w + \varepsilon; x) - f(w; x) \approx$$

$$f(w^{(t+1)}; x) - f(w^{(t)}; x)$$

$$\langle \nabla f(w; x), \varepsilon \rangle.$$

$$f(w^{(t)} - \underbrace{\eta \nabla L(w^{(t)})}_{-\varepsilon}; x) - f(w^{(t)}; x) = -\eta \langle \nabla f(w^{(t)}, x), \underbrace{\nabla L(w^{(t)})}_{\text{ }}$$

Defining  $f(w^{(t)}, x_i) = u_i^{(t)}$ ,

$$= -\eta \left\langle \nabla f(w^{(t)}, x), \sum_{i=1}^N (u_i^{(t)} - y_i) \nabla f(w^{(t)}, x_i) \right\rangle$$

$$= -\eta \sum_{i=1}^N (u_i^{(t)} - y_i) \underbrace{\left\langle \nabla f(w^{(t)}, x), \nabla f(w^{(t)}, x_i) \right\rangle}_{H^{(t)}(x, x_i)}$$

for any  $x$ ,

$$\boxed{f(w^{(t+1)}, x) - f(w^{(t)}, x)} = -\eta \cdot \underbrace{\sum_{i=1}^N (u_i^{(t)} - y_i) \hat{H}^{(t)}(x; x_i)}_{\text{Kernel regression with a time-varying "tangent" Kernel.}}$$

$\Rightarrow$  Moral: gradient descent  $\equiv$  Kernel regression with a time-varying "tangent" Kernel.

→ For very wide NAs, Kernel remains "nearly fixed" with time.   
 ↘ initialized randomly (with Gaussian)



For "large" amount of time,  $w^{(t)}$  doesn't change much!

$w^{(t)}$  can change very slightly & still achieve zero error..

→  $w$ 's don't correspond to "useful" features.

(random features  
[Rahimi; Recht])



## **KERNEL IN THE INFINITE WIDTH LIMIT**