# CS 229br: Advanced Topics in the theory of machine learning

HAR

### Boaz Barak





Ankur Moitra MIT 18.408

Yamini Bansal **Official TF** 

**Dimitris Kalimeris** Unofficial TF

HAKVARD



Gal Kaplun



Preetum Nakkiran Unofficial TF Unofficial TF







### What is deep learning?









God's view:  $f = \operatorname{argmax}_{f} \mathbb{E}[ \bigoplus(f) \mid data ]$   $f = \operatorname{argmax}_{f} \mathbb{E}[ \bigoplus(f \leftrightarrow w) ]$   $w \sim W \mid data$ 



## This seminar

- Taste of research results, questions, experiments, and more
- Goal: Get to state of art research:
  - Background and language
  - Reading and reproducing papers
  - Trying out extensions
- Very experimental and "rough around the edges"
- A lot of learning on your own and from each other
- Hope: Very interactive in lectures and on slack



CS 229br: Survey of very recent research directions, emphasis on experiments



MIT 18.408: Deeper coverage of foundations, emphasis on theory

## Student expectations

- Not set in stone but will include:
- Pre-reading before lectures
- Scribe notes / blog posts (adding proofs and details)
- Applied problem sets (replicating papers or mini versions thereof, exploring) Note: Lectures will not teach practical skills – rely on students to pick up using suggested tutorials, other resources, and each other. Unofficial TFs happy to answer questions!
- Some theoretical problem sets.
- Might have you grade each other's work
- Projects self chosen and directed.

Grading: We'll figure out some grade – hope that's not your loss function 🙂

## Rest of today's lecture:

Blitz through classical learning theory

- Special cases, mostly one dimension
- "Proofs by picture"



PATTERNS, PREDICTIONS, AND ACTIONS

A story about machine learning

Moritz Hardt and Benjamin Recht

## Part I: Convexity & Optimization

## Convexity

- 1. Line between (x, f(x)) and (y, f(y)) above f.
- 2. Tangent line at x (slope f'(x)) below f. 3. f''(x) > 0 for all x



CLAIM: f''(x) < 0 implies tangent line at x above f

Proof:  $f(x + \delta) = f(x) + \delta f'(x) + \delta^2 f''(x) + O(\delta^3)$  implies f below line

## Convexity

- 1. Line between (x, f(x)) and (y, f(y)) above f.
- 2. Tangent line at x (slope f'(x) ) below f.
- 3. f''(x) > 0 for all x



CLAIM: If f above (x, f(x)) - (y, f(y)) then exists z with tangent line at z above f Proof by picture:



### Gradient Descent

 $x_{t+1} = x_t - \eta f'(x_t)$ 

Dimension d:  $f'(x) \rightarrow \nabla f(x) \in \mathbb{R}^d$   $f''(x) \rightarrow H_f(x) = \nabla_2 f(x) \in \mathbb{R}^{d \times d} \text{ (psd)}$  $\eta \sim 1/\lambda_d \text{ drop by } \sim \frac{\lambda_1}{\lambda_d} ||\nabla||^2$ 

Y

$$f(x_t + \delta) \approx f(x_t) + \delta f'(x_t) + \frac{\delta^2}{2} f''(x_t)$$

 $f(x_{t+1}) \approx f(x_t) - \eta f'(x_t)^2 + \frac{\eta^2 f'(x_t)^2}{2} f''(x_t) = f(x_t) - \eta f'(x_t)^2 (1 - \frac{\eta f''(x_t)}{2})$ 

 $x_t x_{t+1}$ 

- If  $\eta < 2/f''(x_t)$  then make progress
- If  $\eta \sim 2/f''(x_t)$  then drop by  $\sim \eta f'(x_t)^2$



 $f(x_{t+1}) \approx f(x_t) - \eta f'(x_t)^2 \left(1 - \frac{\eta f''(x_t)}{2}\right) + \eta^2 \sigma^2 f''(x_t)$ 

- If  $\eta < 2/f''(x_t)$  and (\*)  $\eta \sigma^2 \ll f'(x_t)^2/f''(x)$  then make progress
- If  $\eta \sim 2/f''(x_t)$  and (\*) then drop by  $\sim \eta f'(x_t)^2$

### Part II: Generalization



**Empirical 0-1 loss** 

Generalization gap:  $\mathcal{L}(f) - \hat{\mathcal{L}}_{S}(f)$  for f = A(S)





K (log scale)

## Generalization bounds

 $\log |\mathcal{F}|$  can be replaced with dimensionality / capacity of  $\mathcal{F}$ 

VC dimension: max d s.t.  $\mathcal{F}$  can fit every labels  $y \in \{\pm 1\}^d$  on every samples  $x_1, \dots, x_d$ 

Rademacher complexity: max *d* s.t.  $\mathcal{F}$  can fit random labels  $y \sim \{\pm 1\}^d$  on random  $x_1, \dots, x_d \sim X$ 

PAC Bayes: d = I(A(S); S) where S is training set.

Margin bounds: max *d* s.t. linear classifier satisfies  $\langle w_i, x_i \rangle y_i > ||w_i|| \cdot ||x_i|| / \sqrt{d}$  for correct predictions



### General form: If $C(f) \ll n$ then $\hat{\mathcal{L}}(f) - \mathcal{L}(f) \approx 0$

Intuition: Can't "overfit" – if you do well on *n* samples, must do well on population

Challenge: Many modern deep nets and natural C,  $C(f) \gg n$ Hope: Find better C?

#### UNDERSTANDING DEEP LEARNING REQUIRES RE-THINKING GENERALIZATION

Chiyuan Zhang\* Massachusetts Institute of Technology chiyuan@mit.edu

Benjamin Recht<sup>†</sup> University of California, Berkeley brecht@berkeley.edu Samy Bengio Google Brain bengio@google.com

Moritz Hardt Google Brain mrtz@google.com

Oriol Vinyals Google DeepMind vinyals@google.com

**ICLR 2017** 

### General form: If $C(f) \ll n$ then $\hat{\mathcal{L}}(f) - \mathcal{L}(f) \approx 0$

Intuition: Can't "overfit" – if you do well on *n* samples, must do well on population

Challenge: Many modern deep nets and natural C,  $C(f) \gg n$ Hope: Find better C?

#### Performance on CIFAR-10

Classifier	Train	Test
Wide ResNet <sup>1</sup>	100%	92%
ResNet 18 <sup>2</sup>	100%	87%
Myrtle <sup>3</sup>	100%	86%

#### Performance on noisy CIFAR-10:\*

Train	Test
100%	10%
100%	10%
100%	10%

#### <sup>1</sup> Zagoruyko-Komodakis'16 <sup>2</sup> He-Zhang-Ren-Sun'15 <sup>3</sup> Page '19

#### Generalization gap (noise = 0.00)

End-to-end Supervision





1 1

0.8

1

0.9

Bansal-Kaplun-B, ICLR 2021

### "Double descent"



Belkin, Hsu, Ma, Mandal, PNAS 2019.

Nakkiran-Kaplun-Bansal-Yang-B-Sutskever, ICLR 2020

Data:  $(x_i, f^*(x_i) + N)$  where  $f^*$  degree  $d^*$  polynomial

Algorithm: SGD to minimize  $\sum (f(x_i) - y_i)^2$  where  $f \in \mathcal{F}_d$  = degree d polys



degree  $d = C(\mathcal{F})$ 

Data:  $(x_i, f^*(x_i) + N)$  where  $f^*$  degree  $d^*$  polynomial

Algorithm: SGD to minimize  $\sum (f(x_i) - y_i)^2$  where  $f \in \mathcal{F}_d$  = degree d polys



Data:  $(x_i, f^*(x_i) + N)$  where  $f^*$  degree  $d^*$  polynomial

Algorithm: SGD to minimize  $\sum (f(x_i) - y_i)^2$  where  $f \in \mathcal{F}_d$  = degree d polys





Part III: Approximation and Representation











#### After Fourier Transform



## Fourier Transform

$$\int_0^1 |f - g| < \epsilon$$

Every continuous  $f:[0,1] \rightarrow \mathbb{R}$  can be arbitrarily approximated by g of form

$$g(x) = \sum \alpha_j e^{-2\pi i \beta_j x}$$

i.e.,  $f \approx \text{linear in } \varphi_1(x), \dots, \varphi_N(x) \text{ where } \varphi_j(x) = e^{-2\pi i \beta_j x}$  $\varphi: \mathbb{R} \to \mathbb{R}^N \text{ embedding}$ 

For some natural data, representation  $\varphi$  is "nice" (e.g. sparse, linearly separable,...)

Tasks become simpler after transformation  $x \rightarrow \varphi(x)$ 



 $\varphi(x,y) = (xy, x^2, y^2)$ 

Xavier Bourret Sicotte, https://xavierbourretsicotte.github.io/Kernel\_feature\_map.html

## More approximation

 $ReLU(x) = \max(x, 0)$ 

Every continuous  $f: [0,1] \to \mathbb{R}$  can be arbitrarily approximated by g of form  $g(x) = \sum \alpha_i \operatorname{ReLU}(\beta_i x + \gamma_i)$ 

Proof by picture:



## More approximation

Every continuous  $f: [0,1] \to \mathbb{R}$  can be arbitrarily approximated by g of form  $g(x) = \sum \alpha_i \operatorname{ReLU}(\beta_i x + \gamma_i)$ 

Proof by picture:





### Higher dimension $r(x) = \max(\langle \alpha, x \rangle + \beta, 0)$ , d + 1 parameters



$$J_{(1,0)}(x,y) = I(x)$$



### Higher dimension $r(x) = \max(\langle \alpha, x \rangle + \beta, 0)$ , d + 1 parameters



$$J_{(1,0)}(x,y) = I(x)$$



### Higher dimension $r(x) = \max(\langle \alpha, x \rangle + \beta, 0)$ , d + 1 parameters





## How many ReLU's

Every  $f: [0,1]^d \to \mathbb{R}$  can be approximated as sum of  $r_{\alpha,\beta}(x) = \max(\langle \alpha, x \rangle + \beta)$ 

Can discretize  $\mathbb{R}^{d+1}$  to  $O(1)^d$  points – need at most  $O(1)^d = 2^{O(d)}$  ReLUs

(For formal proofs, see Cybenko '89, Hornik '91, Barron '93; MIT 18.408)

THM: For some f's ,  $2^{c \cdot d}$  ReLU's are necessary

Random  $f: [0,1]^d \rightarrow \mathbb{R}$ 

Choose  $f(x) \in \{0,1\}$  independently at random for  $x \in \{1/4, 3/4\}^d$ 

**Extend linearly** 

**LEMMA:** For every fixed  $g: [0,1]^d \to \mathbb{R}$ ,  $\Pr\left[\int |f - g| < 0.1\right] \le 2^{-\epsilon \cdot 2^d}$ 



THM: For some f's ,  $2^{c \cdot d}$  ReLU's are necessary

Random  $f: [0,1]^d \to \mathbb{R}$ 

Choose  $f(x) \in \{0,1\}$  independently at random for  $x \in \{1/4, 3/4\}^d$ 

**Extend linearly** 



**LEMMA:** For every fixed  $g: [0,1]^d \to \mathbb{R}$ ,  $\Pr\left[\int |f-g| < 0.1\right] \le 2^{-\epsilon \cdot 2^d}$ 

LEMMA  $\Rightarrow$  THEORM: Need  $2^{\epsilon \cdot 2^d}$  distinct g's  $\Rightarrow$  need  $\approx 2^d$  distinct ReLUs

PROOF OF LEMMA: For  $x \in \{1/4, 3/4\}^d$  let  $Z_x = 1$  iff  $\int |g(y) - f(y)| dy > 0.2$  over  $y \in [-1/4, +1/4]^d + x$ 

Each  $Z_x$  is independent with expectation at least 1/2  $\Rightarrow \Pr[\frac{1}{2^d}\sum Z_x < 1/4] \le 2^{-0.1 \cdot 2^d}$   $\frac{1}{2^d}\sum Z_x \ge 1/4 \Rightarrow \int |g - f| \ge 1/20$ 

### Bottom line

For every\* function  $f : \mathbb{R}^d \to \mathbb{R}$  many choices of  $\varphi : \mathbb{R}^d \to \mathbb{R}^N$  s.t.  $f \approx Linear_f \circ \varphi$ 

Want to find  $\varphi$  that is useful for many interesting functions f

useful:

Or f's such that  $\psi \circ f$  interesting for some non-linear  $\psi \colon \mathbb{R} \to \mathbb{R}$ 

- *N* is not too big
- Efficiently compute  $\varphi(x)$  or  $\langle \varphi(x), \varphi(y) \rangle$
- $\langle \varphi(x), \varphi(y) \rangle$  large  $\Leftrightarrow x$  and y are "semantically similar"
- For interesting *f*'s, coefficients of *Linear*<sub>*f*</sub> are *structured*.

### Part IV: Kernels



## Kernel methods (intuitively)

Distance measure K(x, x')

Input:  $(x_1, y_1), ..., (x_n, y_n)$  s.t.  $f^*(x_i) \approx y_i$ 

To approx  $f^*(x)$  output  $y_i$  for  $x_i$  closest to x

or output  $\sum \alpha_i y_i$  with  $\alpha_i$  depending on  $K(x, x_i)$ 

## Hilbert Space

Linear space: v + u,  $c \cdot v$ 

Dot product:  $\langle u, v + c \cdot w \rangle = \langle u, v \rangle + c \langle u, w \rangle$ ,  $\langle u, v \rangle = \langle v, u \rangle$ ,  $\langle v, v \rangle \ge 0$ 

Can solve linear equations of form {  $\langle v_i, x \rangle = b_i$  } knowing only  $\langle v_i, v_j \rangle$ 

Also do least-square minimization  $\min_{x} \sum (\langle v_i, x \rangle - b_i)^2$  knowing only  $\langle v_i, v_j \rangle$ 

**DEF:** Sym matrix  $K \in \mathbb{R}^{n \times n}$  is p.s.d. if  $v^T K v \ge 0$  for all  $v \in \mathbb{R}^n$ Equivalently  $\lambda_1(K), \dots, \lambda_n(K) \ge 0$ 

**CLAIM:** *K* is p.s.d. iff\*  $u^T K v$  is inner product

**PROOF** ( $\Rightarrow$ ):  $u^T Kv = \psi(u) \cdot \psi(v)$  where  $\psi(u_1, ..., u_n) = (\sqrt{\lambda_1}u_1, ..., \sqrt{\lambda_n}u_n)$ , expressed in eigenbasis

### Kernel Methods

Goal: Solve linear equations, least-square minimization, etc. under  $\varphi$ 

**Observation:** Enough to compute  $K(x, x') = \langle \varphi(x), \varphi(x') \rangle$ 

Let  $\hat{x} = \varphi(x)$ , given  $\{\hat{x}_i, y_i\}_{i=1..n}$  want to find  $\hat{w} \in \mathbb{R}^n$  s.t.  $\langle \hat{x}_i, \hat{w} \rangle \approx y_i \forall i$ 

can compute  $\widehat{w} = \sum \alpha_i \ \widehat{x}_i$  using  $K(x_i, x_j)$ 

To compute prediction on new x using  $\hat{w}$ , can compute  $\hat{x} = \sum \beta_i \hat{x}_i$  using  $K(x, x_i)$