

## Lecture 10:

### Computational complexity of Deep Learning

From last lecture, we saw that using ERM estimator, we can learn any Lipschitz low-dim model efficiently.

However, when restricting to estimators that are efficiently learned by SGD, we saw that some low-dim fcts are harder to learn and require more steps/sample (for online SGD)

⇒ There seem to be a computational bottleneck for learning large-flop fcts. Is it fundamental?

This lecture is a very short overview of computational hardness in DL.

- \* Present computational/learning hardness for NNs
- \* Conversely, what can NNs learn efficiently?

E.g. by changing hyperparameters, SGD-trained NNs can emulate

→ kernel methods: any problem efficiently learnable by kernel methods can be learned by SGD-trained NNs (<sup>with right</sup> <sub>hyperparam</sub>)

→ support recovery for small leap low-dim fcts

↳ What other problems can be learned by SGD-trained efficiently?  
Are there "easy problems" on which they fail?

# Computational Hardness of NNs

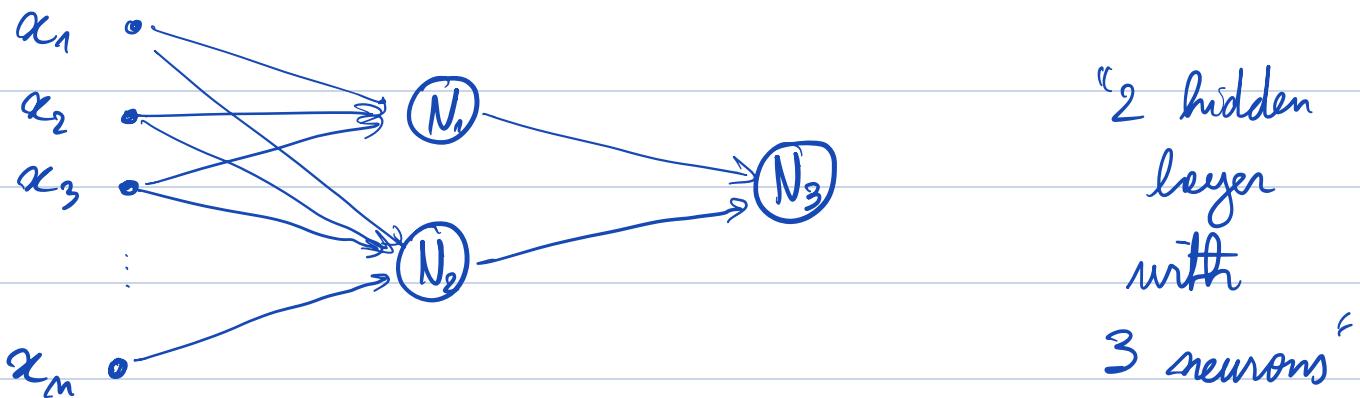
(A) Empirical risk minimization over NNs is NP-complete

Exist result: hardness result for training NNs

→ "we can't hope to have an algo that unconditionally find an (approximate) minimizer of ERM over NNs"  
 (under general hardness conjecture)

Training 3-Node network [Blum, Rivest, 92]

(with many follow-up work)



$$N_i: \varphi_i(z) = \text{sign}(\langle w_i, z \rangle + b_i)$$

(3)

Q\*: Given a set of  $O(d)$  examples  $(x_i, y_i)$  with  $x_i \in \{0,1\}^d$  and  $y_i \in \{-1, 1\}$ . Does there exist a 3-neuron network that has 0 training error?

Thm [Blum, Rivest, '92]: Training 3-neuron NN is NP-complete

This means that if  $P \neq NP$ , no polynomial time algo solve Q\* for any  $(x_i, y_i)$  in  $\text{poly}(d)$  runtime.

Reminder: [Some quick definitions]

"P": problems that can be solved by a deterministic turing machine in polynomial time

"NP": problems that can be solved by a nondeterministic Turing machine (transition fn is not deterministic) in polynomial time

"Polynomial-time reduction": We say that a problem A is as hard as problem B, if there exists a polynomial reduction from B to A. Hence, if I can solve A in polynomial time, then I can solve B in polynomial time. Conversely, if we expect B to be hard to solve, then we expect A to be hard too.

"NP complete": a problem A in NP such that any problem in 4 NP is reducible to A. This means that A is as hard as any problem in NP.

If  $P \neq NP$ , NP-complete problems cannot be solved in polynomial time. ( $P \stackrel{?}{=} NP$  is a major unsolved problem)

Proof idea: We can reduce  $Q^*$  to: given a set of

$\mathcal{O}(d)$  samples does there exist either

- an hyperplane separating  $\pm 1$
- two hyperplanes that split the space in 4 quadrants s.t. one quadrant contains all positive (resp negative) labels and no negative (resp positive) labels

$\tilde{Q}^*$

{

The second case is the "hard" case. We can show that we can reduce a NP complete problem to this problem  
→ making it NP-complete

Set-Splitting: Given  $S$ , a collection of subsets  $\{C_i \mid C_i \subset S\}$

Does there exists  $S_1, S_2$  with  $S_1 \cap S_2 = \emptyset$  and  $S_1 \cup S_2 = S$  such that  $C_i \not\subset S_1$  and  $C_i \not\subset S_2$  for all:

(5)

Thm [ Lovasz (Garey & Johnson, '79)] Set-splitting is  
NP-complete.

We need to construct a reduction from set-splitting to  $\mathcal{Q}^*$

Given an instance of set-splitting:

$$S = \{s_1, \dots, s_d\} \quad C = (C_1, C_2, \dots) \quad C_j \subseteq S$$

we convert it to the following instance of  $\mathcal{Q}^*$

- Let the origin  $(\overbrace{0, 0 \dots 0}^d)$  have label +1

- for each  $s_i$ , we make a point  $x_i = (0, \dots, \overset{i\text{-th}}{1}, \dots, 0)$  and label it  $y_i = -1$

- for all  $C_j = (s_{j1}, \dots, s_{jk})$ , we put a point  $x_{C_j} = (0 \dots 0 \overset{1}{1}, 0 \dots 0 \overset{2}{1} \dots 0 \dots 0 \overset{k}{1} 0)$  and label it  $y_{C_j} = +1$

E.g.  $S = (s_1, s_2, s_3) \quad C_1 = (s_1, s_2) \quad C_2 = (s_2, s_3)$

-1 at	(1 0 0)	(0 1 0)	(0 0 1)
+1 at	(0 0 0)	(1 1 0)	(0 1 1)

6

(Remark: we can always append 0's s.t. #pts is  $O(d)$ )

Prop: This set of pts is separable by a quadrant iff set-splitting has a solution

Proof:  $\Leftarrow$  If set-splitting has  $S_1 \cup S_2 = S, S_1 \cap S_2 = \emptyset$  s.t.  $C_i \notin S_1$  and  $C_i \notin S_2$  for all  $i$ :

Then

consider  $P_1, P_2$  s.t.  $P_j: w_{j1}x_1 + \dots + w_{jd}x_d + \frac{1}{2} \geq 0$

$$\text{s.t. } w_{ji} = \begin{cases} -1 & \text{if } s_i \in S_j \\ d & \text{o.w.} \end{cases}$$

Then:

$$P_j \text{ predicts } \left\{ \begin{array}{l} + \text{ for } (0, \dots, 0) \\ + \text{ for training pts with +} \\ - \text{ for } x_i \text{ if } s_i \in S_j \end{array} \right.$$

Intersection  $P_1 \geq 0, P_2 \geq 0$  (quadrant) contains all the pts with +1 and none of -1 labels.

Let  $S_1$  be the set of points -1 that  
are separated from  $(0, \dots, 0)$  by  $P_1$ .  
Let  $S_2$  be the points -1 not in  $S_1$  that  
are separated from  $(0, \dots, 0)$  by  $P_2$ .

By assumption  $P_1 \geq 0 \wedge P_2 \geq 0$  contain all +1 labels  
and none of -1  
so  $S_1 \cup S_2 = S$  &  $S_1 \cap S_2 = \emptyset$  by construction.

It remains to show that  $C_j \notin S_1, S_2$ .

$C_j = \{s_{j1}, \dots, s_{jk}\}$  and points  $x_{j1}, \dots, x_{jk}$   
label -1

$$x_{C_j} = x_{j1} + \dots + x_{jk} \text{ label } +1$$

If  $C_j \subseteq S_1$  then  $P_1$  separate  $(0, \dots, 0)$  from  $x_{j1}, \dots, x_{jk}$

and in particular it separates  $(0, \dots, 0)$  from  $x_{C_j}$   
contradiction

Same for  $C_j \subseteq S_2$ .

→ Hence if we solve  $\tilde{Q}^*$  we can solve set splitting!



Is this a compelling lower bound?

(should we give up on Deep Learning? :))

For this result:

- \* "worst case" over training data points  
(need to solve given any coll<sup>o</sup>  $(x_i, y_i)$ )
- \* The architecture is fixed: it is NP hard to decide whether there exist weights that perfectly fit the data.

In Deep learning (and machine learning in general), we care instead about: given training data, can we construct a model that predicts well new data point.

In learning theory, this is described as the difference between  
proper vs improper learning

Hypothesis class:  $\mathcal{H} \subseteq \{h: \mathcal{X} \rightarrow \{-1, +1\}\}$

e.g.: the set of all 3-neuron neural networks.

(9)

def: [Proper learning] Given data  $(x_i, y_i) \sim D$  where  
 $y_i = h(x_i)$  with  $h \in H$ , the algorithm output  
 $\hat{h} \in H$  s.t.

$$\text{Test error}_{\mathcal{D}}(\hat{h}) = \mathbb{P}_{\mathcal{D}}(\hat{h}(x) \neq h(x)) \leq \varepsilon$$

def [Improper learning] The algorithm is allowed to output  
 $\hat{h} \notin H$  (e.g. in a much broader class of models)  
s.t.  $\text{Test error}_{\mathcal{D}}(\hat{h}) \leq \varepsilon$

Improper learners are allowed to use a much richer / more expressive hypothesis class even if data comes from a simpler class

Hardness of proper learning does not imply hardness of  
improper learning

Here is an old counterexample:

Def: [k-term DNF] (Disjunctive normal form)

It is the set of boolean formulae

$$\bigvee_{i=1}^k \left( \bigwedge_{j=1}^{m_i} L_{ij} \right)$$

↳ either  $X_j$  or  $\bar{X}_j$   
 $\in \{0, 1\}$

Thm: [Pitt, Valiant] Under a general computational hardness assumption  
 There is no polynomial time algo for learning properly k-term DNFs for  $k \geq 3$

Conversely:

Thm [Valiant] There is a polynomial time algo that learns improperly k-term DNFs.

Here, use k-CNF (conjunctive normal form)

$$\bigwedge_{i=1}^n (L_{i1} \vee L_{i2} \vee \dots \vee L_{ih})$$

We have k-CNF,  $\not\models$  k-DNFs

E.g. 3 DNFs can always be written as 3 CNFs

$$T_1 \vee T_2 \vee T_3 = \bigwedge_{\substack{L_1 \in T_1 \\ L_2 \in T_2 \\ L_3 \in T_3}} (L_1 \vee L_2 \vee L_3)$$

(Convincie yourself)

There are 3 CNFs  
not of this form.

- 3 CNFs can be learned improperly efficiently
  - ↳ algo that always outputs a 3 CNFs that perfectly fits the training data points:
    - \* start with all 3-sized conjunctions ( $\binom{2d}{3}$  of them) and remove all that don't agree with training data
    - \* it generalizes because  $|H_{3-\text{CNF}}| \leq 2^{8d^3}$  and

standard generalization error bound:

$$|\text{error}_{\mathcal{D}}(\hat{h}) - \text{empirical error}(\hat{h})| \leq \varepsilon \text{ w.p } 1-\delta$$

when sample size  $n \geq \frac{1}{2\varepsilon^2} (\ln(2|M|) + \ln(\frac{1}{\delta}))$

[This is "Uniform convergence" bound simply using tail bound + union bound]

Summary: don't trust hardness results for proper learning!

- one can enrich the model class and learn improperly efficiently
- this is what we do in practice with NNs: we take highly overparametrized NNs (much richer than needed to fit the data) which seems to greatly simplify the ERM problem.

Caveat that we saw in lecture 2: by taking a much richer class of functions, we make generalization much harder

- we might require way more samples

There is a fundamental "computational vs statistical trade-off" for proper vs improper learning.

Question: if proper learning hardness is not interesting, can we show hardness of improperly learning NNs?

B

## Hardness of improperly learning NNs

Klivans and Sherstov (2006) showed that for every  $\epsilon > 0$ , it is hard to learn improperly 2-hidden layer neural networks with  $d^\epsilon$  neurons, under a cryptographic hardness assumption.

[This was improved by [Daniely, Linial, Shalev-Shwartz 19] to  $M = \omega_d(1)$  neurons under a different hardness assumption]

What does it mean "cryptographically hard": it's a problem that is not NP-hard, but people have built cryptosystems assuming the problem is hard (e.g. factoring)

Here: shortest vector problem (SVP) in a lattice

$$\text{Lattice} = \{\alpha_1 v_1 + \dots + \alpha_m v_m : \alpha_1, \dots, \alpha_m \in \mathbb{Z}\}$$

$f(n)$ -SVP: approximate the length of the shortest non-zero vector within a factor  $f(n)$ .

Def: [Intersection of  $k$ -halfspaces]  $x \in \{-1, 1\}^d$

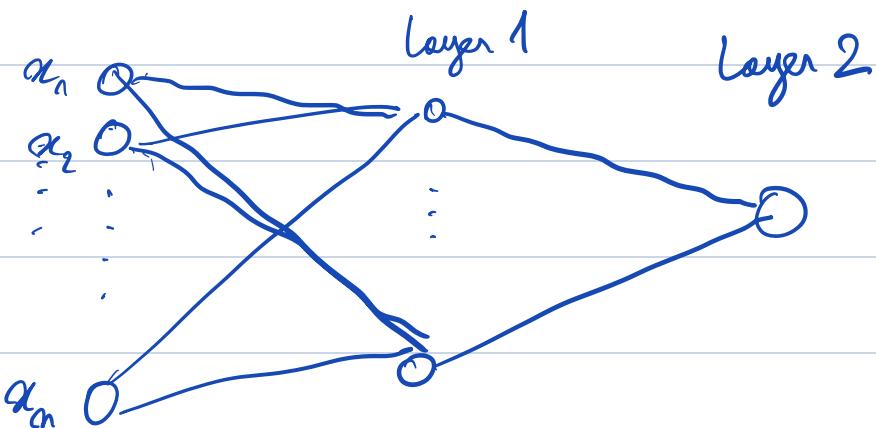
$$\omega_j \in \{-1, 1\}^d$$

$$h(x) = \prod_{i=1}^k \mathbb{1} [\langle \omega_i, x \rangle > 0]$$

$\mathcal{M}_{k\text{-halfspaces}}$  = the set of these fcts

Thm [Klivans, Sherstov, '06] Assuming  $d^{\frac{3}{2}}\text{-SVP}$  is hard,  
then for every  $\varepsilon > 0$ ,  $\mathcal{M}_{d^\varepsilon\text{-halfspaces}}$  cannot be  
learned efficiently, even improperly.

intersection of halfspaces can be written as a 2-hidden-layer  
NN with  $d^\varepsilon$  neurons and ReLU non-linearity



For each half space  $\omega_j \rightarrow 2$  neurons on layer 1

$$\langle \omega_i, x \rangle_+ \quad \text{and} \quad (\langle \omega_i, x \rangle - 1)_+$$

(15)

$$\text{Then } \langle \omega_i, x \rangle_+ - (\langle \omega_i, x \rangle - 1)_+ = \begin{cases} 1 & \text{if } \langle \omega_i, x \rangle > 0 \\ 0 & \text{o.w.} \end{cases}$$

Then in layer 2:

$$\left( \sum_{i=1}^k \langle \omega_i, x \rangle_+ - (\langle \omega_i, x \rangle - 1)_+ + 1 - k \right)_+ \\ = \begin{cases} 1 & \text{iff } x \text{ is in the intersection} \\ 0 & \text{o.w.} \end{cases}$$

Summary: There exists distributions that can be fitted with  $d^\varepsilon$  neurons, that cannot be learned efficiently even if we take much larger neural networks

→ ERM will learn efficiently (approx=0, gen  $\leq \frac{\text{polylog}(d)}{\sqrt{m}}$ )

→ No algo on neural nets will succeed in polynomial time

Is it a compelling hardness result?

→ better than before but still unsatisfying

To prove hardness in this case, they use a weird data distribution. Here is the high level idea:

In public-key encryption system:

take a message  $m \in \{0,1\}^P \rightarrow \text{encr}(m) \in \{0,1\}^q$

↳ public key: use to encrypt message  
→ easy

given  $x = \text{encr}(m)$

$\text{decr}(x) = m$

↳ private key to decrypt the encrypted message

↳ computationally easy if you have private key

hard otherwise

Klivans & Sherstov used a cryptosystem (Regev) based on hardness of SVP. The decoding fact can be implemented as an intersection of  $d^2$  halfspaces.

Idea: If hypothesis class  $H$  contains "decr" function, then it must be hard to improperly learn.

Otherwise an attacker can generate data

$(x = \underline{\text{encr}}(m), y = m = \text{decr}(m))$

public key  
so can compute  $\text{encr}(m)$  efficiently  $\uparrow$  they sample messages

and if he can learn the decryption function even improperly, then it is able to predict the decryption letter more random on a new encrypted message  
 ↳ can break the security of this cryptoscheme.

So... basically the above hardness result is about learning neural nets from data points that are encrypted messages ... not exactly a natural data distribution!!!

Does this hardness hold for more natural data distributions?

- $\text{Unif}(\mathbb{S}^{\pm 13^d})$ : there is  $d^{O(\frac{\log k}{\varepsilon^2})}$  time algo to learn intersection of half-spaces within error  $\varepsilon$ .  
 ↳ in the following, a large class of algo fail to do better.

- Gaussian dist: [Vempala] in  $\text{poly}(d, k, \varepsilon) + k^{O(\frac{\log k}{\varepsilon^2})}$   
 so poly time even when  $k = \frac{\log d}{\log \log d}$   
 [while Amit et al., ∃ dist such that impossible for  $k = w_d(1)$ ]

How to prove computational lower bounds for "nice" distributions?

- it is too much to ask for lower bounds based on NP-hardness or cryptographic hardness
- instead, prove lower bounds for some "large" classes of algorithm

In this lecture, we will consider a powerful framework that has proven successful to capture hardness in many settings

(C)

## Statistical Queries

Kearns (1998) introduced SQ as a framework to design algo that are tolerant to noise.

↳ it restricts algo access to data

In standard ("PAC") learning, an algo has access to  $n$  samples  
 In SQ model, algo only has access to (approximate) oracles on the data, and not directly the samples.

Specifically the algo specifies a query  $\phi: \mathcal{Y} \times \mathcal{X} \rightarrow [-1, 1]$  (19) and a tolerance  $\gamma$  and it gets a response  $r$  s.t.

$$|r - \mathbb{E}_{\mathcal{D}}[\phi(y, x)]| \leq \gamma$$

data dist

⇒ this is a restriction of what algo can do.

Rank: From  $n$  samples, then  $\left| \hat{\mathbb{E}}_n[\phi(y, x)] - \mathbb{E}_{\mathcal{D}}[\phi(y, x)] \right| \leq \frac{1}{\sqrt{n}}$

SQ algo with  $q$  queries of tolerance  $\gamma$ :

- \* proceeds in  $q$  rounds, at each round  $t \in [q]$ , algo issues a query  $\phi_t$  and receive a response

$$|r_t - \mathbb{E}_{\mathcal{D}}[\phi_t]| \leq \gamma$$

$\phi_t$  is chosen adaptively, that is, can depend on past responses

$$v_1, \dots, v_{t-1}$$

- \* After  $q$  rounds, algo output a model.

Rank: any SQ algo with  $q$  queries of tolerance  $\gamma$  can be implemented using  $\mathcal{O}\left(\frac{\log q}{\gamma^2}\right)$  samples

In particular, if a problem can be solved by an SQ algo

- with 1) at most polynomially many queries  
 2) each tolerance at least inverse polynomially big  
 3) can compute next query  $\phi_t$  from  $v_1, \dots, v_{t-1}$   
 and evaluate  $\phi_t$  efficiently

Then there exists a learning algo that succeeds in polynomial time

In summary, SQ algo don't directly look at data but instead use approximate statistics of the data.

⚠ Lower bound: SQ "succeeds" if it succeeds for any compatible

responses  $v_1, \dots, v_q$

↳ hence lower bound is against "worst-case" noise  
 "adversarial"

↳ much more pessimistic than statistical noise

Nonetheless: - Easy to prove lower bounds in this model  
 - Correctly capture lower bound in many settings

⇒ We can prove many tight lower/upper bounds within this framework.

Rank: If a problem is efficiently SQ learnable, then it is efficiently PAC learnable even with a lot of noise in the labels  
 ↳ SQ framework give us a way to design noise-tolerant algo

## Correlation Statistical Queries (CSQ)

A restriction of SQ introduced by Bshouty and Feldman (2002)

$$\phi(y, \alpha) = y \tilde{\phi}(\alpha) \quad \tilde{\phi}: \mathcal{X} \rightarrow [-1, 1]$$

Rank: If  $y \in \{-1\}$ , then any statistical query can be decomposed as

$$\begin{aligned} \phi(y, \alpha) &= \phi(1, \alpha) \frac{1+y}{2} + \phi(-1, \alpha) \frac{1-y}{2} \\ &= \underbrace{\frac{1}{2} y [\phi(1, \alpha) - \phi(-1, \alpha)]}_{\text{CSQ query}} + \underbrace{\phi(1, \alpha) + \phi(-1, \alpha)}_{\text{this does not depend on label}} \end{aligned}$$

Therefore, if we know distribution over  $\alpha$ ,  $SQ = CSQ$  (meaning  $P_X$  fixed)

However  $|Y| > 2$ , then  $SQ$  might be much more powerful than  $CSQ$

[Note: In the distribution free setting  $SQ \gg CSQ$  even if  $|Y|=2$ ]

For gradient descent w.r.t. squared loss

$$\hat{R}_m(\Theta) = \frac{1}{2m} \sum_{i=1}^m (y_i - f(x_i; \Theta))^2$$

GD algo:  $\Theta \in \mathbb{R}^P$

$$\begin{aligned}\Theta^{(k+1)} &= \Theta^{(k)} - \nabla \hat{R}_m(\Theta^{(k)}) \\ &= \Theta^{(k)} + \underbrace{\frac{1}{m} \sum_{i=1}^m y_i \nabla_{\Theta} f(x_i; \Theta^{(k)})}_{\text{vector of } P \text{ CSQ}} - \underbrace{\frac{1}{m} \sum_{i=1}^m f(x_i; \Theta^{(k)}) \nabla_{\Theta} f(x_i; \Theta^{(k)})}_{\text{does not depend on the labels}}\end{aligned}$$

GD with T steps can be implemented by CSQ with  $P^T$  queries and tolerance  $\approx \frac{1}{\sqrt{m}}$



Careful: if reusing the same data, the neural network can implement "sample extraction" and go beyond SQ (more later)

But if online SGD on squared loss, seems to agree with CSQ

→ if loss  $\neq$  squared loss, not CSQ but still SQ

(See later for more discussion about online SGD)

CSQ "captures" GD on squared loss + is very convenient  
to prove lower bounds

Def [Statistical Query dimension] For a class of functions

$$\mathcal{H} \subseteq \{h: \mathcal{X} \rightarrow \mathbb{R}, \|h\|_{L^2} \leq 1\}$$

The SQ-dim of  $\mathcal{H}$  wrt  $P_x$  is the largest integer  $N$  s.t.  $\mathcal{H}$  contains  $N$  functions  $h_j$  s.t.

$$\|h_j\|_{L^2} = 1 \quad \text{and } |\langle h_i, h_j \rangle| \leq \frac{1}{N}$$

High level: CSQ needs  $\tilde{\Omega}(N\epsilon^2)$  queries to learn  $\mathcal{H}$  with accuracy  $\epsilon$ .  
Let's prove a slightly more general result:

Thm: Let  $\mathcal{H}$  be a class of functions that contains  $N$   $h_j$  s.t.

$\|h_j\|_{L^2} = 1$  and  $|\langle h_i, h_j \rangle| \leq \epsilon$ . Then any CSQ with tolerance  $\epsilon$ , number of queries  $q \leq N \frac{\epsilon^2 - \epsilon}{4}$  and  $N \geq \frac{2}{\epsilon}$  has

output of CSQ algo when true func is  $h$ .

$$\sup_{h \in \mathcal{H}} \|\hat{h} - h\|_{L^2} \geq 1 - 2\epsilon$$

Proof: Idea: with  $q \leq N \frac{\varepsilon^2 - \varepsilon}{\varepsilon}$ , we are not able to distinguish between all  $h_j$ .

That is, each query allow to rule out some  $h_j$ 's whose correlation is not  $\varepsilon$ -close. After  $q$  queries, we are not able to rule out a large fraction of them.

Specifically, we will show that after  $q \leq N \frac{\varepsilon^2 - \varepsilon}{\varepsilon}$  there are at least  $k \geq \frac{N}{2}$  fits compatible with the answers.

$$\begin{aligned} \sup_{h \in H} \|\hat{h} - h\|_2^2 &\geq \min_{\hat{h}} \max_{i \in [k]} \|\hat{h} - h_i\|_2^2 \\ &\geq \min_{\hat{h}} \frac{1}{k} \sum_{i=1}^k \|\hat{h} - h_i\|_2^2 \\ &= \frac{1}{k} \sum_{i=1}^k \|h_i - \frac{1}{k} \sum_{j=1}^k h_j\|_2^2 \\ &= \frac{1}{k} \sum_{i=1}^k \|h_i\|_2^2 - \frac{1}{k^2} \sum_{i=j}^k \langle h_i, h_j \rangle \end{aligned}$$

$$= 1 - \frac{1}{k} - \frac{1}{k^2} \sum_{i \neq j} \langle h_i, h_j \rangle$$

$$\geq 1 - \frac{1}{k} - \varepsilon \geq 1 - \frac{2}{N} - \varepsilon$$

$$\geq 1 - 2\varepsilon$$

Now let's show  $\geq \frac{N}{2}$  facts are compatible with the responses.

Noise is adversarial, so it is enough to consider a deterministic sequence of answers that are compatible with  $\geq \frac{N}{2}$  facts

Consider  $r_1 = \dots = r_q = 0$

and the associated sequence of queries  $\phi_1, \dots, \phi_q$   
 (fully determined by the responses)

For each  $t \in [q]$ , define

$$S_t^+ = \{i \in [N] : \mathbb{E}[h_i(x) \phi_t(x)] > z\}$$

$$S_t^- = \{i \in [N] : \mathbb{E}[h_i(x) \phi_t(x)] < -z\}$$

$S_t^+ \cup S_t^-$  are all the  $h_j$  that we are able to rule out at step  $t$

$$N_t^+ = |S_t^+| \quad N_t^- = |S_t^-|$$

By definition,

$$\begin{aligned} N_t^+ z &\leq \sum_{j \in S_t^+} \langle h_j, \phi_t \rangle \stackrel{\leq 1}{\leq} \|\phi_t\|_{L^2} \left\| \sum_{j \in S_t^+} h_j \right\|_{L^2} \\ &\leq \overbrace{\sqrt{\sum_{j \in S_t^+} \|h_j\|_{L^2}^2}} + \sum_{i \neq j \in S_t^+} \langle h_j, h_i \rangle \end{aligned}$$

So that  $(N_t^+ z)^2 \leq N_t^+ + (N_t^+)^2 \varepsilon$

$$\Rightarrow N_t^+ \leq \frac{1}{z^2 - \varepsilon}$$

Some results holds for  $N_E^- \leq \frac{1}{z^2 - \varepsilon}$  so  
that

$$\sum_{t=1}^q N_t^- \leq \frac{2q}{z^2 - \varepsilon}$$

Hence if  $q \leq N \frac{z^2 - \varepsilon}{\varepsilon}$ , then we rule out at most  $\leq \frac{N}{2}$  fits  $\square$

### Applications

From above, to give lower bound, it is sufficient to find a family of nearly orthogonal functions.

(I)

# Hardness of learning parities

$$\mathcal{H} = \{ h(x) = \chi_S(x) : S \subseteq [d] \}$$

$$\chi_S(x) = \prod_{i \in S} x_i$$

$\mathcal{H}$  contains  $2^d$  orthogonal fcts. They are all binary valued so that  $\text{CSQ} = \text{SQ}$

Applying the above general lower bound, we get

Thm: [Blum, Furst, Kearns, Mansour, ...]

Any SQ algo for learning parities must make at least  $2^{\Omega(d)}$  queries or have  $\mathcal{Z} = 2^{-\Omega(d)}$

Is this lower bound correctly capturing hardness of learning parities?

(1)

There is a polynomial time algo for learning parities: it suffices to solve a linear system

$$m \in \mathbb{F}_2 = \{\pm 1\} \quad A c = y$$

$$\in \mathbb{R}^{n \times d} \quad \begin{bmatrix} -x_1- \\ \vdots \\ -x_n- \end{bmatrix} \quad c = 1_s$$

But this algo is highly non robust and will fail if there is noise in the label

$\neq$  SQ algo that is robust to noise.

② SQ correctly captures lower bound for noisy problem?  
 → not quite, but almost?

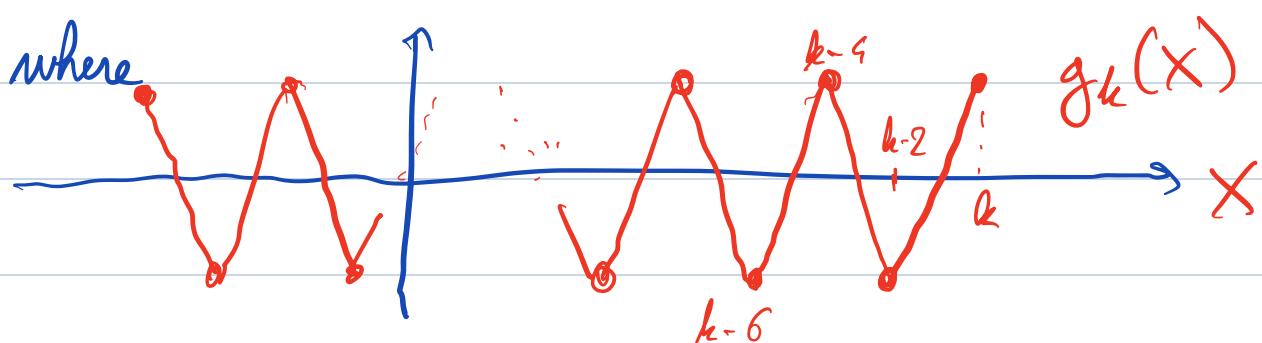
[Blum, Kothari, Wasserman] There is a  $2^{\Theta(\frac{d}{\log d})}$  time algo for learning noisy parities.  
 (→ noise: flip label with prob  $\frac{1}{3}$ )

→ this is slightly better than  $2^{\Omega(d)}$  in the lower bound

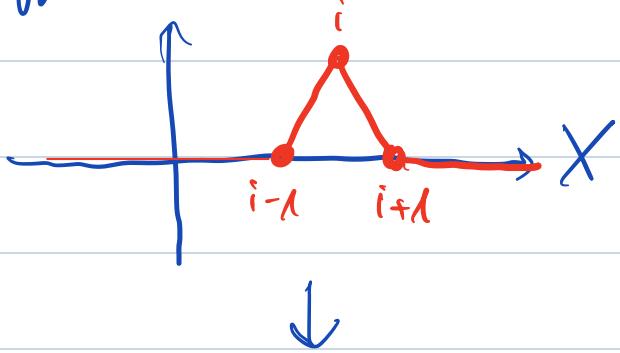
Lower bound for NNs:

Lemma: Any parity function  $\chi_S(x)$   $|S|=k$  can be easily implemented as a 2-layer NN with ReLU activation and  $\Theta(k)$  neurons.

Proof:  $\chi_S(x) = \prod_{i \in S} x_i = g_k\left(\sum_{i \in S} x_i\right)$



It suffices to show that we can implement



→ and then add and subtract these facts

$$(x - i+1)_+ - 2(x - i)_+ + (x - i-1)_+$$

↳ hence taking  $x = \langle 1_S, \alpha \rangle$ , and summing them indeed implement the parity

□

Hence • ERM learns parity fcts:

$$\text{approx error} = 0$$

$$\text{gen error} \leq \frac{\text{poly}(d)}{\sqrt{n}}$$

weights bounded by  $d$   
and at most  $\mathcal{O}(d)$  many

- However GID fail (if we trust we cannot efficiently learn noisy parities)

## II Learn complexity

The set  $\{X_S(C_n) : |S|=k\}$  has  $\binom{d}{k} = \mathcal{O}(d^k)$  orthogonal fct

$$\text{SQ lower bound: } \frac{q}{Z^2} \geq c d^k$$

For more general sparse fct:

$$g(y) = \underbrace{y_1}_{\text{red}} + \underbrace{y_1 y_2 y_3 y_4}_{\text{green}}$$

$$\|g\|_{L^2}^2 = 2$$

If we want to learn  $\|\hat{f} - g\|_{L^2} < 1$ , we need to learn  $y_2 y_3 y_4$ . Conditional on knowing  $y_1 =$  learning 3-parity

and we can prove CSQ lower bound

(31)

$$g_{Z^2} \geq c d^3$$

actually can be made  
SQ lower bound  
for this specific  
fct.

More generally:  $k$ -loop fits,

CSQ requires  $\frac{g}{Z^2} \geq c d^k$

For online SGD with batch 1, if we burst this lower bound, that means we need

$$Td \geq d^k \quad \text{step size of batch 1}$$

$$(z \approx 1)$$

that is  $T \geq d^{k-1}$

It agrees with numerical simulations!!

III

Ridge functions

Take  $\alpha \sim N(0, \text{Id})$  and consider the set

$$\mathcal{H} = \left\{ h(\omega) = \text{He}_k(\langle \omega_0, \omega \rangle) : \omega \in \mathbb{S}^{d-1} \right\}$$

where  $\text{He}_k$  is the degree- $k$  Hermite polynomial  
(normalized)

[recall  $\mathbb{E}_G[\text{He}_k(G) \text{He}_j(G)] = \mathbb{I}_{[j=k]} \quad G \sim N(0, 1)$ ]

Let's construct nearly orthogonal fits with  $\{\omega_j\}_{j=1}^N$

Note that  $\langle h_i, h_j \rangle = \mathbb{E}[\text{He}_k(\langle \omega_i, \omega \rangle) \text{He}_k(\langle \omega_j, \omega \rangle)]$

$$= \langle \omega_i, \omega_j \rangle^k$$

Take a family s.t.  $|\langle \omega_i, \omega_j \rangle| \leq \delta$  for all  $i \neq j$

$\hookrightarrow$  we can choose  $N \geq \exp(c'd\delta^2)$  of them

Then setting  $\varepsilon = \delta^k$  and using general lower bound from above,

CSQ requires  $q \geq \frac{e^{cd\delta^2}}{4} (z^2 - \delta^k)$

Fixing  $\delta = \left(\frac{z^2}{2}\right)^{\frac{1}{k}} \rightarrow \log(8q/z^2) \geq c'dz^{\frac{1}{k}}$

so we need  $\frac{1}{z^2} \gtrsim d^{k/2}$  for  $\frac{q}{z^2}$  to not 33

be exponential in  $d$ . If we heuristically set

$$Z = \frac{1}{\sqrt{m}}$$

$m = \# \text{ of samples}$  (to implement CSQ)  
algo

we need

$$m \gtrsim d^{\frac{k}{2}}$$

samples in order to be efficiently  
solvable

Rank:



This lower bound is incorrect for  $M_k$   
if we reuse samples / learn with another loss  
function. In fact  $m = \Theta(d)$  (info-theoretic)  
is sufficient.

↳ however, it seems to be correct for some functions  
whose first non-zero hermite coeff is  $k$ .

[Dominici et al., 24]

## Noisy Gradient Descent and CSQ

Population GD + gaussian noise

For  $t=1, \dots, T$

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \nabla_{\mathbf{w}} R(\mathbf{w}^t) + \eta \mathbf{z}^t$$

where  $R(\mathbf{w}) = \mathbb{E}[(y - f(x, \mathbf{w}))^2]$

$\sim_{iid} N(0, \sigma^2 I_p)$

Rank: [Online SGD]

$$\mathbf{z}^t = \nabla_{\mathbf{w}} (y^t - f(x^t, \mathbf{w}^t))^2 - \nabla_{\mathbf{w}} R(\mathbf{w}^t)$$

mean-0 iid noise but not Gaussian

Noisy GD can be approximated by a CSQ algo

$$-\nabla_{\mathbf{w}} R(\mathbf{w}) = \underbrace{2\mathbb{E}[y \nabla_{\mathbf{w}} f(x, \mathbf{w})]}_{P \text{ CSQ queries}} - \underbrace{2\mathbb{E}[f(x; \mathbf{w}) \nabla_{\mathbf{w}} f(x, \mathbf{w})]}_{\text{marginal over } x \text{ is fixed}}$$

Let's couple noisy GD with CSQ with tolerance  $\epsilon$

$$\tilde{m}^{k+1} = \tilde{m}^k - \gamma n^k + \gamma \tilde{z}^t \quad \text{some Gaussian noise}$$

$$\hookrightarrow \|n^k - \nabla_{\tilde{\Theta}} R(\tilde{\Theta}^k)\|_\infty \leq \gamma$$

We compare the law of  $\tilde{m}^0, \tilde{m}^1, \dots, \tilde{m}^T \equiv \tilde{m}^{\leq T}$

to the law of  $\tilde{m}^0, \tilde{m}^1, \dots, \tilde{m}^T \equiv \tilde{m}^{\leq T}$

We consider the KL divergence between the two

$$D_{KL}(P \parallel Q) = \int_{s \in S} \log\left(\frac{dP}{dQ}(s)\right) dP(s)$$

We have

$$D_{KL}(\tilde{m}^{\leq T} \parallel \tilde{m}^{\leq T}) = D_{KL}(\tilde{m}^0 \parallel \tilde{m}^0) \quad (= 0 \text{ (some initialized)})$$

$$\text{chain rule for KL} \quad \left( + \sum_{t=1}^T D_{KL}(\tilde{m}^t | \tilde{m}^{\leq t-1} \parallel \tilde{m}^t | \tilde{m}^{\leq t-1}) \right)$$

$$= \sum_{t=1}^T \underbrace{D_{KL}(\tilde{m}^t | \tilde{m}^{t-1} \parallel \tilde{m}^t | \tilde{m}^{t-1})}_{\equiv C_t} \quad ) \text{Markov property}$$

Using the update equation and def of conditional KL

$$C_t = \mathbb{E}_{\omega \sim h(\Theta^{t-1})} \left[ D_{KL} \left( \underbrace{\mathcal{D}_{\Theta}(\Theta - \eta \nabla R(\Theta) + \eta z)}_{N(\Theta - \eta \nabla R(\Theta), \eta^2 \sigma^2 I)} \| \underbrace{\Theta - \eta n + \eta z}_{N(\Theta - \eta n, \eta^2 \sigma^2 I)} \right) \right]$$

$$D_{KL}(N(\mu, \sigma^2 I) \| N(\mu', \sigma^2 I)) = \frac{1}{2\sigma^2} \|\mu - \mu'\|^2$$

$$= \frac{1}{2\sigma^2} \|\nabla R(\Theta) - n\|_2^2$$

$$\leq \frac{1}{2\sigma^2} \|\nabla z\|^2$$

Therefore:  $D_{KL}(\Theta^{\leq T} \| \widehat{\Theta}^{\leq T}) \leq \frac{T \|\nabla z\|^2}{2\sigma^2}$

"total variation distance"

Consider  $\overline{TV}(P \| Q) = \sup_A |P(A) - Q(A)|$

$$\begin{aligned} TV(\Theta^T \| \widehat{\Theta}^T) &\leq TV(\Theta^{\leq T} \| \widehat{\Theta}^{\leq T}) && \text{data processing inequality} \\ &\leq \sqrt{\frac{1}{2} D_{KL}(\Theta^{\leq T}, \widehat{\Theta}^{\leq T})} && \text{ Pinsker's inequality } \end{aligned}$$

$$\leq \frac{\sqrt{PT}}{2} \frac{z}{\sigma}$$

(37)

We consider the optimal coupling between  $\mathbb{W}^T$  and  $\tilde{\mathbb{W}}^T$

$$\sup_{\text{joint prob}} P(\mathbb{W}^T = \tilde{\mathbb{W}}^T) = 1 - TV(\mathbb{W}^T \| \tilde{\mathbb{W}}^T) \geq 1 - \frac{\sqrt{P^T}}{2} \frac{z}{6}$$

Take  $z = \frac{2\delta\sigma}{\sqrt{P^T}}$ :

with probability  $\geq 1 - \delta$ , we have  $\mathbb{W}^T = \tilde{\mathbb{W}}^T$

Hence, if  $\tilde{\mathbb{W}}^T$  fails with  $T_p$  queries of tolerance  $\frac{2\delta\sigma}{\sqrt{P^T}}$

Then, so does noisy GD with probe  $1 - \delta$

Corollary: Noisy GD on polynomial sized NN must take either  $e^{O(d)}$  many steps or have noise  $\sigma = e^{-O(d)}$  in order to learn parities.

Remark: Note that it is still not a lower bound directly on GD or SGD. However still very informative!  
(it is hard to go beyond SQ with noise-tolerant algo)

Rmk: Similar argument: instead of coupling noisy GD to a CSQ algorithm, one can couple to a dynamics with label  $y$  independent of  $x$  (a "junk flow" where there is nothing to learn)  
 → the guarantees obtained are very similar to the ones obtained with CSQ  
 [Abbe, Sandon, '20] [Abbe, Boix-Adsera, '22]

Rmk: Can we prove unconditional lower bounds for SGD directly?  
 Results below indicates that no (if we don't constrain hyperparameters) as it would imply  $P \neq NP$

①

## Poly-time invertibility of DL

Did we learn anything about deep learning from the previous lower bounds?

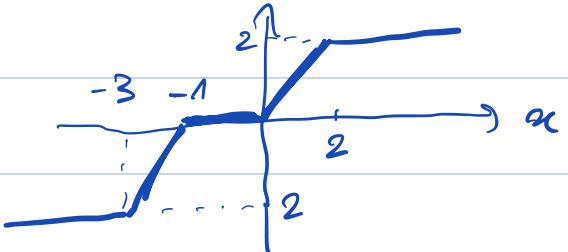
So far, we only used that there exist some problems that are hard to learn with any algo (or any SQ algo) and therefore SGD-trained NNs cannot succeed too..

Instead, we would like to get some insights on why neural networks specifically are so successful and what are their limitations

e.g. is there a problem that is efficiently solvable in polynomial time but where SGD on NNs fail?

Albe and Sordon (2020)  
answered this question negatively in the following sense.

For the following activation function



such that the following holds.

Thm [Abbe, Sondor (2020)] If a problem can be learned in polynomial time, then there exists a neural network with the step-wise ramp activation and of poly size such that running SGD on it learn the problem in poly time.

SGD on NNs is P-complete

Proof idea: it is an emulation proof. One can implement "bit extraction" that allows the NN to extract samples using SGD and then poly-computable function can be implemented as a polynomial-sized circuit.

The proof is quite involved and not that insightful.

E.g. SGD on NNs can learn noiseless properties in polytime by emulating Gaussian elimination.

Rank: For the above emulation, one needs enough "precision" at each step. If the precision is too small, SGD on NNs can only implement poly-time SQ algo.

E.g. will fail to learn pointies.

Is it an interesting result?

Yes: proving an unconditional lower bound on SGD-trained NNs would imply  $P \neq NP$ .

→ so we cannot show such a result without solving one of the major open problems of our time

No: These emulation NNs are highly non-standard and require complex architectures and initialization  
 ↳ these hyperparameters depend on the problem.

↳ This is very different than what NNs do in practice  
 (i.e., they are not emulating Gaussian elimination...)

Instead, NNs used in practice have regular architectures and generic initializations

What we really want: understand what SGD on regular NNs can learn / fail to learn.

Some work in that direction (e.g. learning leap fit).

The challenge of regular NNs vs emulsion NNs is that it requires to understand SGD dynamics from generic initialization on regular architectures...

→ high-dimensional non-convex dynamics

This is hard!

Next lecture: high-dimensional non-convex dynamics