

Lecture 1 [Room will change!]

Plan for today:

- * Mathematical setup
- * Course schedule & logistics
- * Further setup, first theorems?

In practice, a deep network a set of functions \mathcal{F} , s.t.

- | | |
|---|--|
| <p>[standard.]</p> <ol style="list-style-type: none"> ① $\mathcal{F}: \mathcal{X} \rightarrow \mathcal{Y}; w \in \mathbb{R}^D$
for fixed parameter dim d. ② \mathcal{F} has uniformly bounded runtime (in terms of elementary ops on real numbers) | <p>[less standard]</p> <ol style="list-style-type: none"> ③ Efficient hardware ④ convenient programming libraries ⑤ Amenable to GD. |
|---|--|

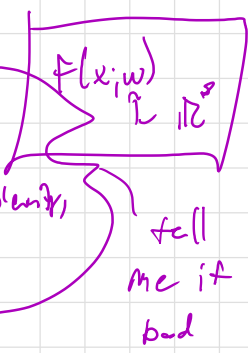
$\hat{R} =$ "perf on seen data" | $R =$ "perf on unseen data"

Examples

- Ⓐ DALL-E input: english sentence | outputs 1024x1024x3
 \mathbb{R}^{256} | # data $\mathbb{R} \approx 2^{30}$
- Ⓑ AlphaFold input: amino acid sequence \rightarrow 3-d structure
 \hat{R} & R hard to obtain, but well-defined but "no common distribution"

Scope of this class: Not just simpler models, but specific focus:

- ① "feedforward networks", typically with 2-layers
- ② Goal is to show $R(\hat{f})$ small via a specific error decomposition
output of alg
unseen performance
- ③ Short proofs & reasonable facts
- ④ Avoid looseness via adaptive complexity
- ⑤ Bridge old & new.



Defn. (Feed forward network)

$$x \mapsto \sigma_2(W_2 \sigma_1(W_1 x + b_1) + b_2) \dots$$

weights bases
matrix bias
nonlinearity, transfer, activation

Csp. $x \mapsto a^T \sigma(Wx)$ (2-layer, no bias)

or $\sigma_i: \mathbb{R}^{m_i} \times \mathbb{R}^{m_{i+1}}$ are coordinate-wise scalar functions ($m_i = m_{i+1}$)

ReLU $z \mapsto \max\{0, z\}$, sigmoid $z \mapsto \frac{1}{1+e^{-z}}$

"Error decomposition"

Want: $R(\hat{f})$ small; all we know is $\hat{R}(\hat{f})$ smallish
 helper: $\bar{f} \in \mathcal{F}$ "nearly best over R "

$$R(\hat{f}) = \underbrace{R(\hat{f}) - \hat{R}(\hat{f})}_{\text{generalization}} + \underbrace{\hat{R}(\hat{f}) - \hat{R}(\bar{f})}_{\text{optimization}} + \underbrace{\hat{R}(\bar{f}) - R(\bar{f})}_{\text{generalization}} + R(\bar{f})$$

"approximation"

Remark "interpolation" $\hat{R}(\hat{f}) = 0 \ll R(\bar{f}) \approx R(\hat{f})$

- * web page has everything
- * places
 - * in-person; OHI outside (TR, 30min each)
 - * zoom available, de-emphasized
 - * gradescope
 - * discussion edstem ("scheduled")

- * Grading
 - * 80% in 4 homeworks, typed, submit in gradescope, don't cheat

- * proj
 - * course notes
-

Apk
 shallow / constructive apk
 initialization & overparameterization
 non-shallow construction

Opt
 NTR smooth / strong convexity } 6 lectures
 (mean-field?)
 non-NTR margin analysis

[??
 , ,
 generalization?
 other topics