

Ch 20 Neural Networks: approximation

Friday, March 27, 2020 3:05 PM

Neural Networks or "Artificial Neural Network" (ANN)

Originally inspired by brain, but now we don't need the physical link
i.e., \neq Hodgkin-Huxley models...

We'll focus on feed-forward ANN, and only cover basic theory

Not covering Pooling, batch-normalization, dropout, recurrent ANN, auto-encoders, generative-adversarial networks (GAN), ...

Feed-forward ANN is described by its architecture (a DAG (a directed acyclic graph w , nodes V , edges E)
and an activation function σ \curvearrowright node = neuron
and the edge weights w . \curvearrowright hyperparameters (train via grad. student descent)
 \curvearrowright parameters (train via SGD)

$$w: E \rightarrow \mathbb{R}, w \in \mathcal{H}_{V, E, \sigma}$$

Assume a layered structure, $V = \bigcup_{t=0}^T V_t$, $|V_0| = n+1$ is input layer \curvearrowright for constant

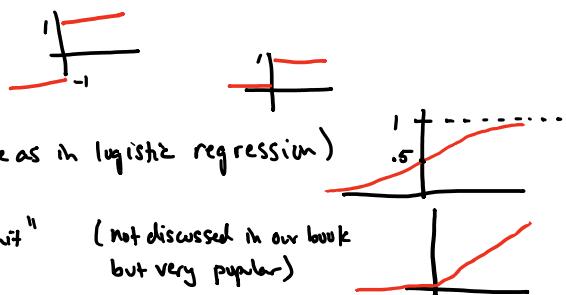
$\sigma: \mathbb{R} \rightarrow \mathbb{R}$, sometimes write $\sigma(\vec{a})$ to mean $(\sigma(a_i))_i = \sigma(a_i)$, i.e., apply coordinate-wise
typical examples:

Ex 1a: $\sigma(a) = \text{sign}(a)$ "sign" fn

Ex 1b: $\sigma(a) = \mathbb{1}_{[a>0]}$ "threshold" fn

Ex 2: $\sigma(a) = \frac{1}{1+e^{-a}}$ "sigmoid" fn (Same as in logistic regression)

Ex 3: $\sigma(a) = \text{ReLU}_+ = \max(0, a)$ "rectified linear unit" (not discussed in our book
but very popular)



many more variants, but similar theoretical properties

(but nonlinearity usually essential)

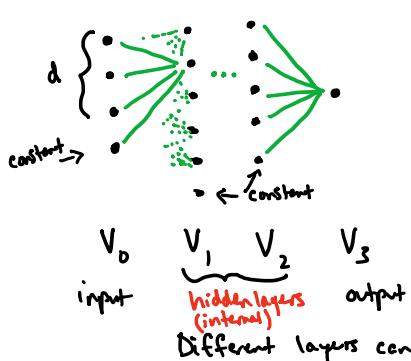
To calculate the output, let V_{ti} be i^{th} neuron in t^{th} layer,

which has input a_{ti} and gives output $o_{ti} = \begin{cases} \sigma(a_{ti}) & i \neq 0 \\ 1 & i=0 \end{cases}$ (special "constant" neuron)

where input is $a_{ti} = \sum_r W((V_{t-1r}, V_{ti})) \cdot o_{t-1r}$

(note: all input/output really depend on original input vector x)

Much simpler w/ a picture



$$T = \# \text{layers (excluding input)}, \text{ e.g. } 3$$

$$\text{width} = \max_t |V_t|, \text{ e.g., } 6$$

$$\text{size} = |V|$$

output layer is 1 neuron for binary classification

but more for k-way classification (a major benefit of ANN is avoiding the ugly reductions of **multiclass** to binary via, e.g., one-vs.-all or all-pairs, cf. ch 17)

or d for autoencoder, $x \in \mathbb{R}^d$

Different layers can have different sizes

Fully connected means complete bipartite for adjacent layers
(vs a **convolutional** layer, for example)

Outline of our neural net unit:

I. Approximation Error

II. Estimation Error (generalization)

III. Optimization Error (new!)

Let $A(S)$ be output of neural net w/ training data S like ERM

$$L_D(A(S)) = \underbrace{\min_{\text{h} \in \mathcal{H}_{\text{NN}}} L_D(h)}_{\text{(I) Approx.}} + \underbrace{\min_{\text{h} \in \mathcal{H}_{\text{NN}}} \hat{L}_S(h)}_{\text{(II) Estimation}} - \underbrace{\min_{\text{h} \in \mathcal{H}_{\text{NN}}} L_D(h) + L_D(A(S)) - \hat{L}_S(A(S))}_{\text{also (II) Estimation}} + \underbrace{\hat{L}_S(A(S)) - \min_{\text{h} \in \mathcal{H}_{\text{NN}}} \hat{L}_S(h)}_{\text{(III) Optimization}} \quad (= 0 \text{ if ERM})$$

this is "classical". More "modern" approach incorporates training algorithm into analysis

(before, algorithm can only hurt you by giving large optimization error,
now, algo. can help by improving generalization (e.g., stability analysis) and you
don't have to include estimation error)

- Optimization Error hard to bound
- Estimation error of large ANN is huge if you do, e.g., VC-dim analysis

"Lack of theoretical understanding of ANN" is due to 1) these pts.

Neural nets generalize better than expected, and the type of algorithm (ex SGD w/ dropout)
makes a difference

MVHT new theory since ~2015: one of biggest open problems

Basically, theory predicts a finite sample complexity, but it's too large / conservative

History:

Ideas, algorithms, analysis 1980's

Graeff Hinton, Yann LeCun, Yoshua Bengio (Mike Mozer at CU)

- but not too much practical success (not bad, but not clear winner)
 until ~2012,
 (1) conv nets showed excellent performance
 (2) "deep" (many layers), enabled by more computational power
 (3) GPUs, better computers
 (4) tricks (dropout, ...)
 (5) success \Rightarrow more success (more researchers, more tricks/insight ...)
 (6) software (tensorflow, caffe, pyTorch, etc.)

① The expressive power of ANN (i.e., approximation error)

Lemma For approximation, more-is-better: expressive power is nondecreasing as we add more layers or depth.

Proof intuitive. If we add depth, can always set weights = 0

If we add layers, for $\sigma(a) = \mathbb{1}_{[a>0]}$, note $\sigma^2(a) = \sigma(a)$, so make weights the identity

Approximation theorems

"Approximate all $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ " is ill-posed

(1) need a metric (eg, L^2 , ptwise)

(2) metric restricts type of functions f

(eg., need $\int f^2(x)dx < \infty$, or need f defined ptwise)

Classical thm: Weierstrass Approx. thm, says $P([a,b])$ dense in $C([a,b], \|\cdot\|_\infty)$
 (or Stone-W.)
 ↑ all polynomials

For binary classification, $f: \mathbb{R}^n \rightarrow \{\pm 1\}$, ϵ -approx. wrt $\|\cdot\|_\infty$ doesn't make sense,
 but exact representation also undesirable (since uncountable cardinality)

Better: exact representation, but for $f: \{\pm 1\}^n \rightarrow \{\pm 1\}$ (Boolean Fcn.)

(and if $f: \mathbb{R}^n \rightarrow \mathbb{R}$, represent \mathbb{R} w/ 64-bits on computer, then take $n=64d$)

Claim 20.1

$(\forall n \in \mathbb{N}) (\exists V, E \text{ of depth 2, width } 2^n + 1) \text{ s.t. } \mathcal{H}_{V, E, f=\text{sign}}$ contains all Boolean fns $f: \{\pm 1\}^n \rightarrow \{\pm 1\}$

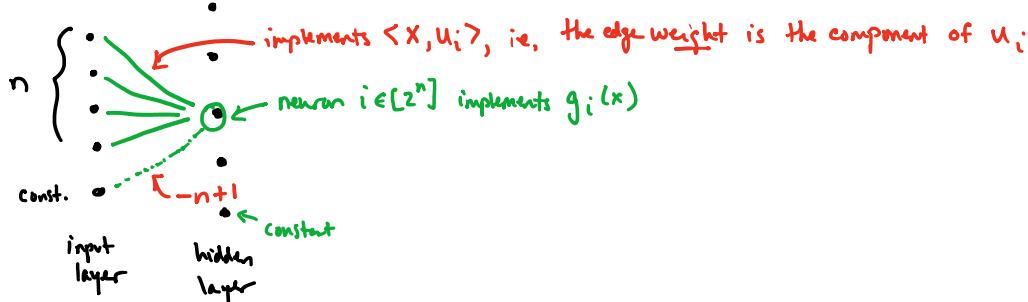
Proof

Enumerate $\{\pm 1\}^n = \{u_i\}_{i \in [2^n]}$. Observe $V_i, \forall x \in \{\pm 1\}^n, \langle x, u_i \rangle = \begin{cases} n & x=u_i \\ n' \leq n-2 & x \neq u_i \end{cases}$

Since $x \neq u_i \Rightarrow$ differ in at least one coordinate, so that "+1" is "-1", a difference of 2

for all $i \in [2^n]$, define

$$g_i(x) = \text{sign}(\langle x, u_i \rangle - n + 1) = \begin{cases} +1 & x=u_i \\ -1 & x \neq u_i \end{cases}$$



Now, let $f: \{\pm 1\}^n \rightarrow \{\pm 1\}$ be arbitrary. It is characterized by the inputs that give an output of +1. Let $\{u_i\}_{i \in F}$ be these inputs, $|F| \equiv k \leq n$

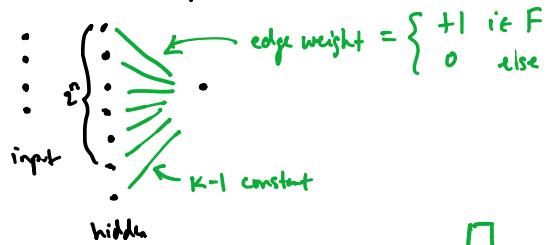
$$\text{Observe } f(x) = \begin{cases} +1 & \text{if } x \in \cup_{i \in F} u_i \\ -1 & \text{else} \end{cases}$$

$$= \text{sign} \left(\sum_{i \in F} g_i(x) + k - 1 \right)$$

$$\text{Since } x \in \cup_{i \in F} u_i \Rightarrow \sum_{i \in F} g_i(x) = 1 - (k-1) = 2-k \Rightarrow 2-k+k-1 > 0$$

$$x \notin \cup_{i \in F} u_i \Rightarrow \sum_{i \in F} g_i(x) = -k \Rightarrow -k+k-1 < 0$$

So for our last layer we implement



□

What about real-valued functions?

Example result

Exercise 20.1 Let $f: [-1, 1]^n \rightarrow [-1, 1]$ be ρ -Lipschitz. Then ($\forall \varepsilon > 0$)

\exists ANN f_{NN} w/ $\sigma = \text{sigmoid}$ st. $\|f - f_{NN}\|_\infty < \varepsilon$

proof sketch

wlog let f be ρ -Lipschitz wrt $\|\cdot\|_\infty$ norm on \mathbb{R}^n
(ie., adjust value of ρ if needed)

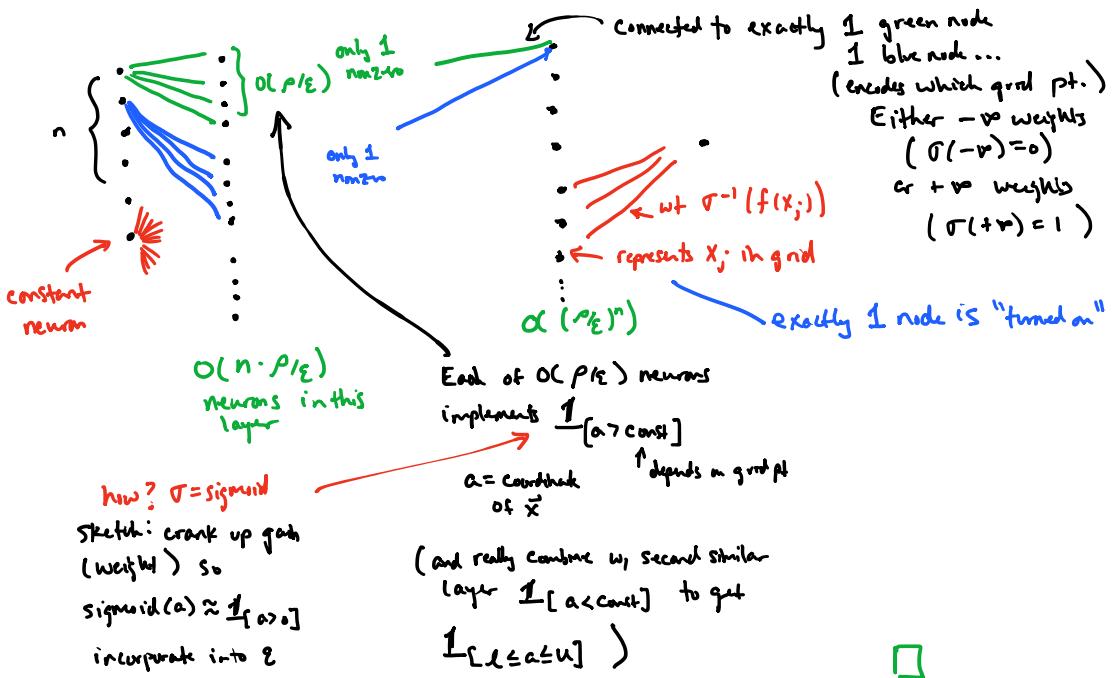
Make grid on $[-1, 1]^n$, spacing of $O(\varepsilon/\rho)$ in each dimension

(so $O(\rho/\varepsilon)$ grid pts/dimension, or $O((\rho/\varepsilon)^n)$ grid pts overall)

hence $x, x' \in$ same grid cell

$$\Rightarrow |f(x) - f(x')| \leq \rho \|x - x'\|_\infty \leq \varepsilon/2$$

Now, given $x \in [-1, 1]^n$, can we determine which grid cell it belongs to?



Back to Claim 20.1, the width is $2^n + 1$ (not grid!)

Not just due to proof technique

Theorem 20.2 ($\forall n \in \mathbb{N}$) if $H_{V,E,\sigma=\text{sign}}$ contains all Boolean
 fn $f: \{0,1\}^n \rightarrow \{0,1\}$ (or $\{\pm 1\}$, book's type) then
 $|V|$ is exponential in n . Similar results if $\sigma = \text{sigmoid}$

Proof

Let $H_{V,E,\sigma=\text{sign}}$ contain all such Boolean functions, thus it can shatter

$m = 2^n$ vectors $\Rightarrow V\text{cdim}(H) \geq 2^n$

As we'll show later, $V\text{cdim}(H) \leq O(|E| \log(|E|))$
 $\leq O(|V|^3)$

so $2^n \leq O(|V|^3) \Rightarrow |V| \geq \Omega(2^{n/3})$ □

Similarly for \mathbb{R} -valued functions:

Thm 20.5 $\forall \epsilon \in (0,1)$, let $H_{V,E,\sigma=\text{sigmoid}}$ be large enough to approximate
 all 1-Lipschitz functions $f: [-1,1]^n \rightarrow [-1,1]$ to within ϵ . Then

$|V|$ is exponential in n .

Related theorems in the literature

Zhou Lu et al. "The Expressive Power of Neural Networks: A View from the Width" (NIPS 2017)

Thm (density in $L^1(\mathbb{R}^n)$) Let $f \in L^1(\mathbb{R}^n)$ then $\forall \varepsilon > 0 \exists \text{ANN } f_{NN}: \mathbb{R}^n \rightarrow \mathbb{R}$
 (fully connected, ReLU activation, depth $4n+1$, width $n+4$) s.t. $\|f - f_{NN}\|_{L^1} < \varepsilon$

G. Cybenko. "Approximation by superpositions of a sigmoidal function" (1989).
 DOI: 10.1007/BF02551274

Thm (density in $C(\bar{\Omega}, \| \cdot \|_\infty)$) Let $\bar{\Omega}$ be any compact subset of \mathbb{R}^n ,
 then ANN w/ one hidden layer, arbitrary depth, any continuous sigmoidal activation
 are dense in $C(\bar{\Omega}, \| \cdot \|_\infty)$

(Issue: f_{NN} w/ σ sigmoidal is bounded, but let σ at last layer be identity)

$\lim_{t \rightarrow +\infty} \sigma(t) = 1$
 $\lim_{t \rightarrow -\infty} \sigma(t) = 0$

Dmitry Yarotsky. "Optimal approximation of continuous functions by very deep ReLU networks" (2018) arXiv: 1802.03620

Thm (density in $C^k(\bar{\Omega}, \| \cdot \|_\infty)$, quantitative) ANN w/ $\sigma = \text{ReLU}$ is dense in $C^k(\bar{\Omega}, \| \cdot \|_\infty)$
 and for error ε , need $O(\log(\frac{1}{\varepsilon}))$ layers and $O((\frac{1}{\varepsilon})^{k/n} \log(\frac{1}{\varepsilon}))$ depth. $|E|$ weights
 in \mathbb{R}^n

David Rolnick and Max Tegmark. "The power of deeper networks for expressing natural functions" (2017) arXiv: 1705.05502

Thm lowerbounds For approximation of \mathbb{R}^n function, ANN needs
 width $\approx \begin{cases} 2^n & 1 \text{ hidden layer} \\ n^{1/k} & k \text{ hidden layers } (1 < k < \text{cutoff}) \\ n & \text{many hidden layers} \end{cases}$

Chulhee Yun, Suvrit Sra, and Ali Jadbabaie. "Small ReLU networks are powerful memorizers: a tight analysis of memorization" (2018) arXiv: 1810.07770

Instead of learning functions, can talk about memorizing datasets

= good, from approx. perspective = bad, from generalization perspective

Thm Shallow A 3-layer ReLU network w/ $\Omega(\sqrt{m})$ hidden nodes can perfectly memorize dataset
 $(O(\sqrt{m}))$ is necessary and sufficient \rightarrow (2 hidden layers of \sqrt{m} each, w/ m points or, product is m)

Back to our book

We saw some Boolean functions require an ANN of exponential size.

Can some be expressed in ANN of polynomial size? Yes.

Thm 20.3 Let $T: \mathbb{N} \rightarrow \mathbb{N}$. $\exists c, b \in \mathbb{R}_+$ s.t. $(\forall n \in \mathbb{N}) \exists (V, E)$ of size $|V| \leq c \cdot T(n)^k + b$ s.t. any function f which can be implemented on a Turing machine in runtime $T(n)$ or less is inside the class $\mathcal{H}_{V, E, r=\text{sign}}$

No proof, but key idea is that an ANN w/ sign can implement the basics of a circuit i.e., a set of universal gates, like $\{\text{NAND}\}, \{\text{NOR}\}, \{\text{AND}, \text{NOT}\}$

In particular, ANN can recreate

Conjunction \wedge	(AND)
disjunction \vee	(OR)
negation \neg	(NOT)

Lemma 20.4

Let $\sigma = \text{sign}$, and let a neuron V have K incoming edges with values $o_i, i \in [K]$. Then (allowing for a constant input neuron) we can adjust weights w_i s.t.

$$\text{the output } \sigma(\sum_i w_i o_i + c) = \bigwedge_i o_i \quad \text{or} \quad \bigvee_i o_i$$

Proof $\bigwedge_{i \in [K]} o_i = \text{sign}(1 - K + \sum_{i=1}^K o_i)$ and $\bigvee_{i \in [K]} o_i = \text{sign}(K - 1 + \sum_i o_i)$

Geometric Intuition

$f: \mathbb{R}^2 \rightarrow \{\pm 1\}$ w/ $\sigma = \text{sign}$. What can we represent?

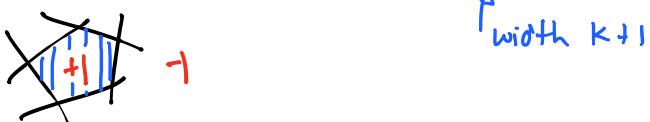
Depth 2 (1 input layer, 1 hidden layer, 1 output layer)

Every neuron in hidden layer is of form $\text{sign}(\langle w, x \rangle + b)$

So it's a halfspace predictor. Then the output can be a conjunction or disjunction of these, so intersection/union of halfspaces.

(If $|V_1| = k+1$, intersection of k halfspaces

↑
for constant i.e., convex polytope w/ k faces

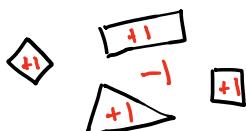


↑
width $k+1$

Depth 3

Take depth 2 and also do conjunctions/disjunctions

So... union of polytopes (or intersection... so nonconvex polytopes)



So you can see it's quite expressive compared to the classes we've studied before (half-spaces, SVM, decision stumps, decision trees)