

THEORY OF MACHINE LEARNING

LECTURE 19

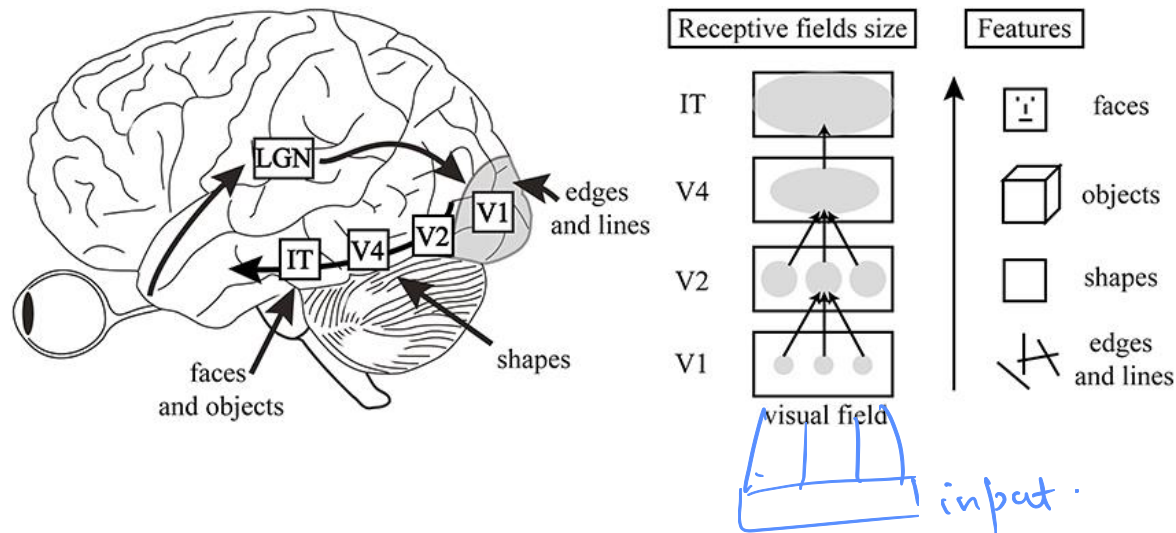
NEURAL NETWORKS – REPRESENTATION, OPTIMIZATION

groups of 2

- HW3 will be out this weekend.
- Project proposal → find a paper you think is "worth presenting" as project, write one para abt it.
- No lecture this Thursday..

RECAP

- Idea behind neural nets:
 - Perceptrons detect “basic” or “primitive” features; ‘composing’ them allows for complex decision-making
 - Supported by human visual system (V1, V2, ...)



RECAP: ARTIFICIAL/DEEP NEURAL NETWORK (DNN)

- **Definition.** A layered "circuit" that takes a vector of input features x , produces output $y = F_r \circ F_{r-1} \circ \dots \circ F_1(x)$, where each F_i is a function of the form $F_i(z) = \sigma(Az + b)$, for some activation function $\sigma()$ (that acts coordinate-wise)

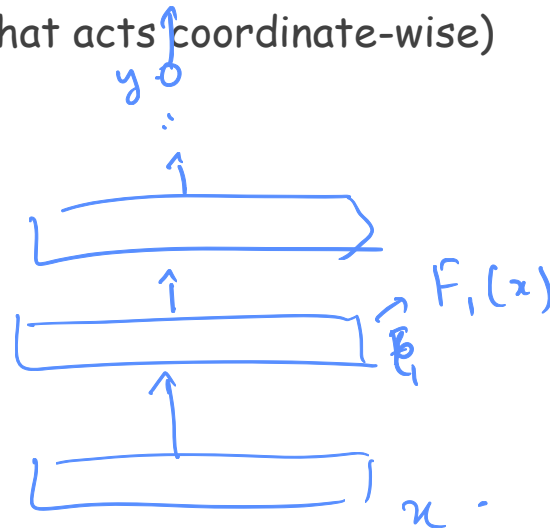
- Common activation functions:

- Threshold

- Sigmoid: (continuous approx.) $\frac{1}{1+e^{-x}}$

- ReLU, Tanh

- ...



LEARNING NEURAL NETWORKS

- Defines a hypothesis class
- **Question (vanilla supervised learning):** given data $(x_1, y_1), (x_2, y_2), \dots$ from some distribution D , find h in this class that minimizes the risk
- ERM problem usually called neural network “training” - given data, find best hypothesis $(f(x_i) = y_i)$ for all i

THEORY OF DEEP LEARNING – THREE BROAD DIRECTIONS

- Expressibility
 - What kinds of functions can be obtained using a DNN? *deep neural net.*
- Training complexity & training dynamics for GD and variants
 - Can the ERM problem be solved efficiently? What guarantees are possible?
- Generalization
 - What kind of generalization bounds can we prove? (VC dimension?) *(complexity of functions)*
no focus on training
If you get m training examples, can you ensure that test error (assuming same dist) is $< f(m)$
 $\frac{1}{\sqrt{m}}$?
params \approx # data samples..

Key: "easy" answers for all questions, but unsatisfactory for realistic settings

EXPRESSIBILITY BASICS

- **Barron's theorem [93]**. Any continuous function f that satisfies an appropriate "niceness" condition (parametrized by C) can be approximated to error ϵ (in $L2$!) by a 2-layer NN with $\sim \frac{C^2}{\epsilon}$ internal nodes
- (Nice functions can be approximated by small NNs)
- **Universal approximation** [Cybenko, Hornik '87,'91]. Any continuous function (over a compact domain) can be approximated by a 2-layer NN with any non-linearity (not a polynomial)

Curse of dimensionality for Cybenko (not Barron)

Moral:

depth-2 nets are universal (express ANY function — given large enough width.)

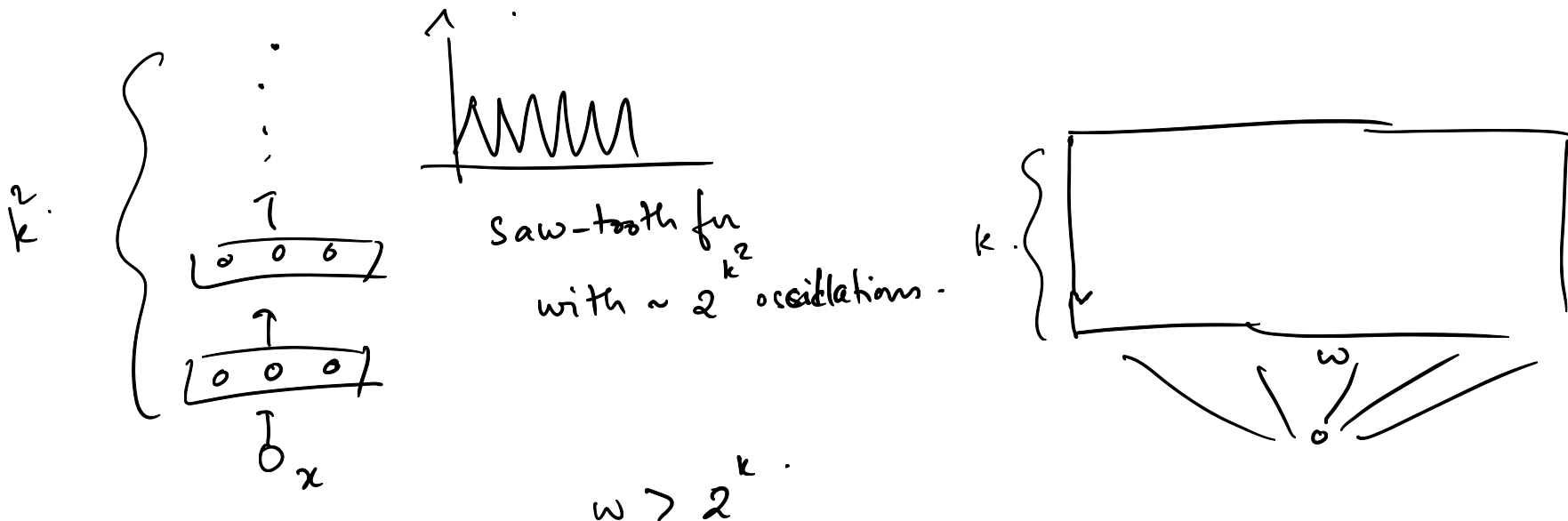
WHY “DEEP” NETWORKS?

- Practical intuition:
 - *high-level.* Depth allows “meaningful features” while width is for “brute force memorization”
- Universality results degrade rapidly with dimensions
 - Curse of dimensionality
 - Modern nets work with high dimensional data
- Does higher depth lead to higher expressibility (with much fewer neurons)?
- Yes! [Eldan and Shamir, Telgarsky]

POWER OF DEPTH

Theorem [Telgarsky 16]. There exists a network of depth k^2 and $O(1)$ width that computes function f , with the property that any network of depth k that approximates f requires width $> 2^k$

(For more general piecewise poly functions, first bound changes to k^3)



(1) A function with $\sim 2^{k^2}$ oscillations can be output by a depth k^2 net with 3 neurons / layer.

that is the output of a ReLU

(2) Any function ~~with~~ width w , ~~depth~~ depth k network is composed of $\leq (2w)^k$ piecewise linear "pieces".

(3) ~~If~~ ^{posl} A function that ~~is~~ with fewer than $2^{k/2}$ pieces cannot even approximate f_n from (1) to an error $\sim \frac{1}{4}$.

$$(2w)^k > 2^{k/2} \Rightarrow w \gtrsim \frac{2^{k/2}}{10}.$$

[Similar results expected for other fun classes.]

MORALS

$$C : \int |w| |\hat{f}(w)| dx$$

(Baron's thing)

■ Depth allows capturing "complex patterns" (oscillations) .
different behaviors in

■ Width allows capturing "different regions of space"

■ What is the right network for an application?

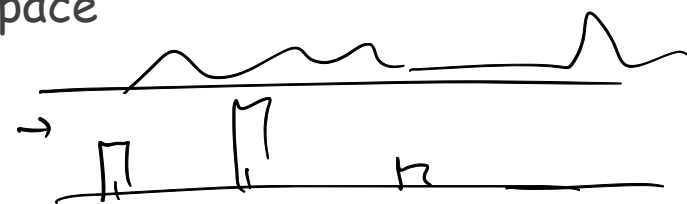
■ Very hard question (Neural Architecture Search)

■ Example of Vision + NLP problems

2013 ■ Hebbian principle → 1950s. → "Neurons that fire together wire together"

■ Needs exploiting domain knowledge (physics informed ML)

[Arora, Ma, Risteski]: deeper networks have a higher "Baron ^{constant} ~~coefficient~~".
(grows exp(d))



- Mode collapse.

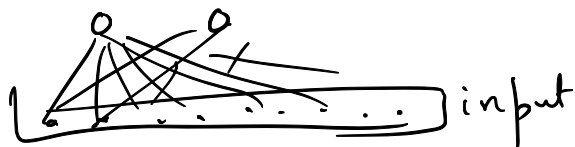
→ Circuit complexity.

NEURAL NETWORK TRAINING

- Supervised learning of NN:** given data $(x_1, y_1), (x_2, y_2), \dots$ from some distribution D , find h that minimizes the ~~empirical~~ risk $(\text{empirical is what we care about.})$.

input \nearrow
 label \nearrow
 (x_0, y_0)
 - Standard metrics: squared loss, cross entropy ($h_w(x)$ is the output of NN).
 find w , so as to minimize $\sum_{i=1}^n (h_w(x_i) - y_i)^2$ $w \rightarrow$ weights.

$$l(h_w(x_i), y_i)$$
 - ERM problem for neural nets
 - NP-hard to learn weights, even if classification error is 0 and we just have 3 internal nodes
- [For the squared loss, ^{ERM} problem is NP hard, because h_w is non-linear.]



COMMON ALGORITHM – GRADIENT DESCENT

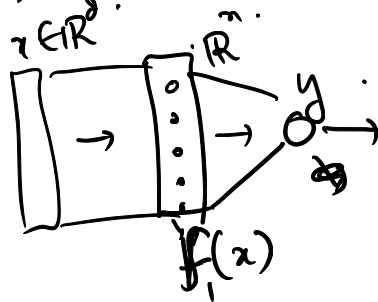
→ Recall how we solved ERM problems:

$$(x_1, y_1), \dots, (x_n, y_n).$$

$h(w; x)$: hypothesis ~~was~~ defined by "weight vector" w on ~~the~~ input x .

(linear hypothesis, $h(w; x) = \text{dot} \langle x, w \rangle + c$)

NN?



$$F_1(x) = \sigma(Ax + b)$$

$$y = \sigma(v^T F_1(x) + b') \quad (\text{model parameters}).$$

weight "vector" w consists of

$$w = (A, v, b, b')$$

In general, weights \equiv (tuple of matrices, biases, ...)

ERM: find w so that $\sum_{i=1}^n \ell(h(w; x_i), y_i)$ is min

given training data, this is $G(w)$; ~~also~~

→ If G were convex, gradient descent (GD) provably works...

→ @ NNs: G is usually (for most choices of σ) not convex.

→ GD still finds "local opt".

Hope: if with sufficient data, finds a "good" local opt

★ "Quality" of "local opt" → diff. between local & global... } "Landscape analysis".

★ # of local opt: →

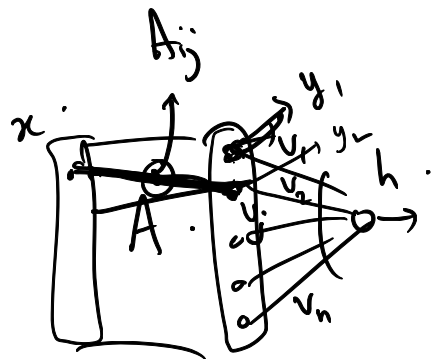
★ Lots of global opt (or "nearly" global opt), but some generalize way better than others

How to do gradient: $w^{(0)}$

$$w^{(t+1)} = w^{(t)} - \eta \nabla G(w^{(t)})$$

$$G = \sum_{i=1}^N (h(w; x_i) - y_i)^2$$

For simplicity, assume $G = (h(w; x) - y)^2$



$$\nabla G = 2(h(w; x) - y) \nabla h(w; x)$$

$$h(w; x) = \sigma(v^T (\sigma(Ax + b)) + b')$$

$$= \sigma(v_1 y_1 + v_2 y_2 + \dots + v_n y_n + b')$$

∇h has components along each entry of A, b, b', v .

how does h change if you vary w ?

$\frac{\partial h}{\partial v_i}$ is "essentially" y_i (if σ were ReLU).

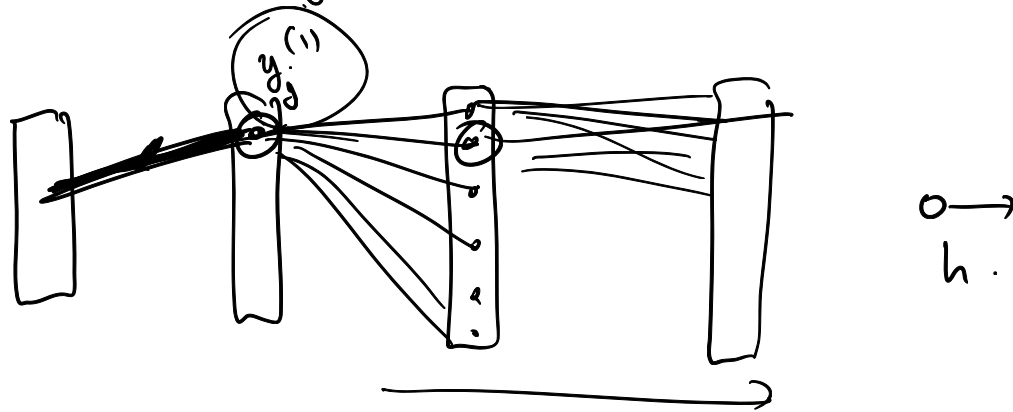
$\frac{\partial h}{\partial v_i} \rightarrow$ one parameter.

$$\sigma' \left(\underbrace{v^T (\dots)} \right) \cdot y_i$$

How do you diff. wrt. A_{ij} ?

$$\frac{\partial h}{\partial A_{ij}} = \sigma' \left(v^T (\dots) \right) \cdot \frac{\partial y_j}{\partial A_{ij}} v_j$$

* If you have > 2 layers



[Hinton et al. '86]

Backpropagation algorithm \rightarrow allows you to compute gradients in \sim linear time.

[Note on Back-propagation.]