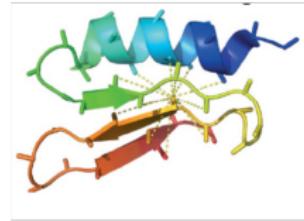
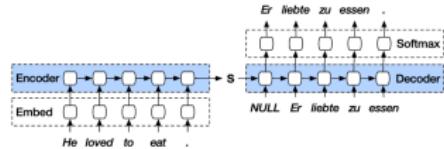
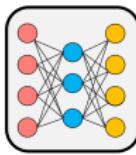


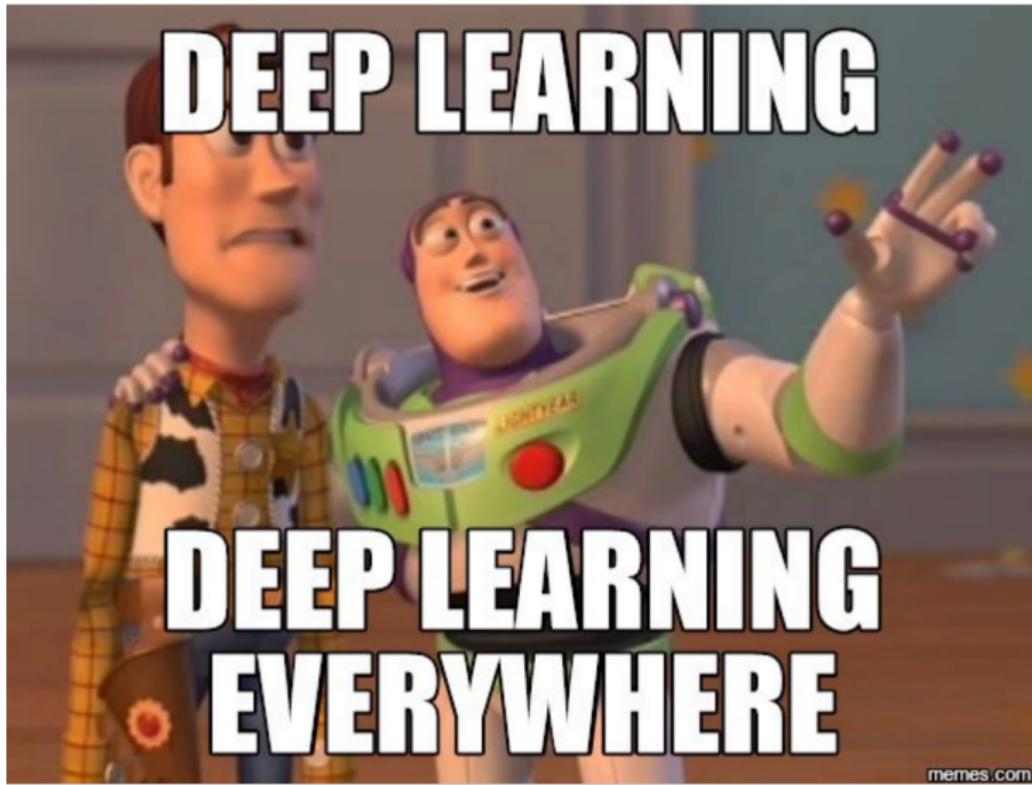
CS 446/ECE 449: Machine Learning

Shenlong Wang

University of Illinois at Urbana-Champaign, 2024

Deep Learning





memes.com

Overview of the next five lectures

- Single layer perceptron
- Multi-layer perceptron / Backpropagation
- Convolutional neural network
- Sequential modeling
- Transformers

Goals of this lecture

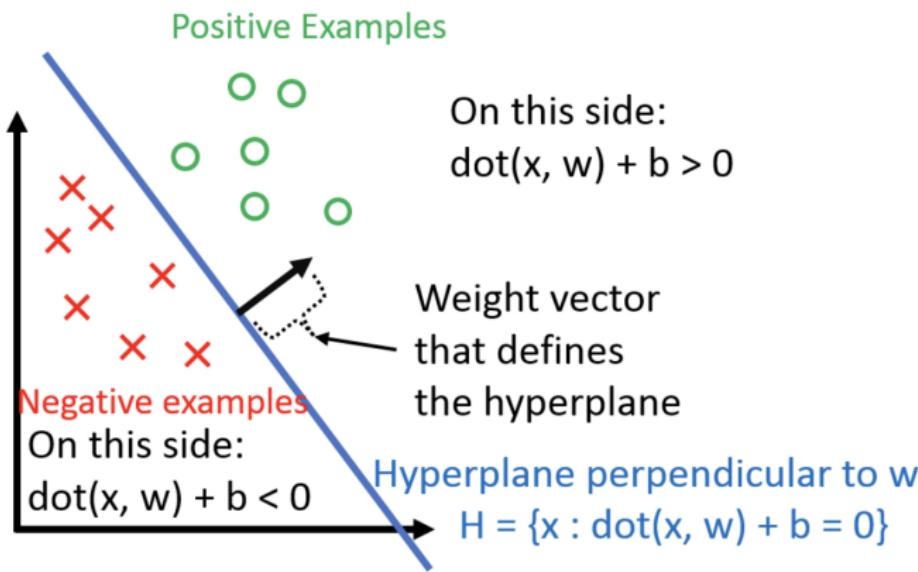
- Getting to know perceptron
- Getting to know deep nets
- Understanding forward pass
- Learning about deep net components

Reading material

- I. Goodfellow et al.; Deep Learning; Chapters 6-9

Perceptron

$$y = \text{sign}(\mathbf{w}^T \mathbf{x}) \text{ let's use } \mathbf{x} \leftarrow [\mathbf{x}, 1]$$



Note that

$$y_i(\mathbf{w}^T \mathbf{x}_i) > 0 \iff \mathbf{x}_i \text{ is classified correctly}$$

How can we learn?

Perceptron algorithm (Frank Rosenblatt, 1958):

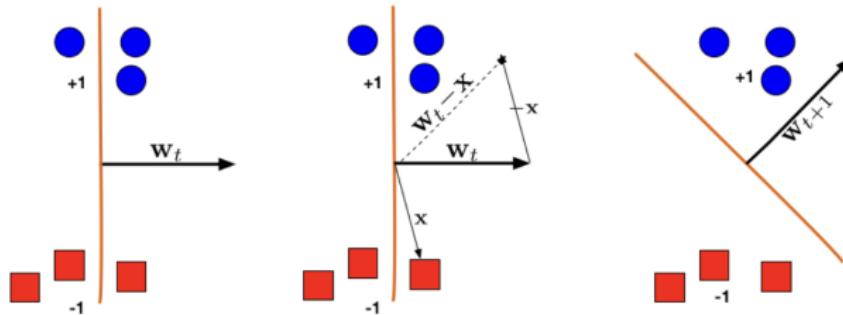
$$y = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

Perceptron Algorithm

```
Initialize  $\vec{w} = \vec{0}$                                 // Initialize  $\vec{w}$ .  $\vec{w} = \vec{0}$  misclassifies everything.  
while TRUE do                                     // Keep looping  
     $m = 0$                                          // Count the number of misclassifications,  $m$   
    for  $(x_i, y_i) \in D$  do                         // Loop over each (data, label) pair in the dataset,  $D$   
        if  $y_i(\vec{w}^T \cdot \vec{x}_i) \leq 0$  then       // If the pair  $(\vec{x}_i, y_i)$  is misclassified  
             $\vec{w} \leftarrow \vec{w} + y_i \vec{x}$              // Update the weight vector  $\vec{w}$   
             $m \leftarrow m + 1$                          // Counter the number of misclassification  
        end if  
    end for  
    if  $m = 0$  then                                 // If the most recent  $\vec{w}$  gave 0 misclassifications  
        break                                       // Break out of the while-loop  
    end if  
end while                                         // Otherwise, keep looping!
```

What does a perceptron update do?

If $y(\mathbf{w}^T \mathbf{x}) \leq 0$, then $\mathbf{w}_{\text{new}} \leftarrow \mathbf{w} + y\mathbf{x}$

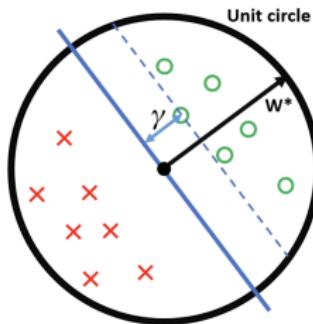


$$\begin{aligned}y(\mathbf{w}_{\text{new}}^T \mathbf{x}) &= y(\mathbf{w} + y\mathbf{x})^T \mathbf{x} \\&= y(\mathbf{w}^T \mathbf{x}) + \| \mathbf{x} \|^2 \\&> y(\mathbf{w}^T \mathbf{x})\end{aligned}$$

Quiz: Assume a data set consists only of a single data point $\{(\mathbf{x}, +1)\}$. How often can a perceptron misclassify this point \mathbf{x} repeatedly?

$$y(\mathbf{w} + k\mathbf{y}\mathbf{x})^T \mathbf{x} \leq 0 \rightarrow k \leq -y(\mathbf{w}^T \mathbf{x})/\|\mathbf{x}\|^2$$

Perceptron Convergence



Assumption:

- All inputs \mathbf{x}_i live within the unit sphere. (after rescaling)
- There exists a separating hyperplane defined by \mathbf{w}^* , with $\|\mathbf{w}\|^* = 1$ (i.e. \mathbf{w}^* lies exactly on the unit sphere.)
- γ is linear separation margin.

Theorem: If all of the above holds, the perceptron algorithm makes at most $1/\gamma^2$ updates.

Perceptron Convergence

Theorem: If all of the above holds, the perceptron algorithm makes at most $1/\gamma^2$ updates.

Proof: Consider the effect of an update:

If $y(\mathbf{w}^T \mathbf{x}) \leq 0$, then $\mathbf{w}_{\text{new}} \leftarrow \mathbf{w} + y\mathbf{x}$, we have:

- $y(\mathbf{x}^T \mathbf{w}) \leq 0$: \mathbf{x} is misclassified by \mathbf{w} - otherwise wouldn't update.
- $y(\mathbf{x}^T \mathbf{w}^*) \geq 0$: \mathbf{w}^* is a separating hyper-plane

(i) Consider the effect of an update on $\mathbf{w}^T \mathbf{w}^*$:

- $(\mathbf{w} + y\mathbf{x})^T \mathbf{w}^* = \mathbf{w}^T \mathbf{w}^* + y(\mathbf{x}^T \mathbf{w}^*) \geq \mathbf{w}^T \mathbf{w}^* + \gamma$
- This means: $\mathbf{w}^T \mathbf{w}^*$ grows by at least γ .

(ii) Consider the effect of an update on $\mathbf{w}^T \mathbf{w}$:

- $(\mathbf{w} + y\mathbf{x})^T (\mathbf{w} + y\mathbf{x}) = \mathbf{w}^T \mathbf{w} + \underbrace{2y(\mathbf{w}^T \mathbf{x})}_{<0} + \underbrace{y^2(\mathbf{x}^T \mathbf{x})}_{0 \leq \leq 1} \leq \mathbf{w}^T \mathbf{w} + 1$
- This means: $\mathbf{w}^T \mathbf{w}$ grows by at most 1.

Perceptron Convergence

Now recall that we initialize $\mathbf{w} = \mathbf{0}$. Hence, initially $\mathbf{w}^\top \mathbf{w} = 0$ and $\mathbf{w}^\top \mathbf{w}^* = 0$ and combining (i) and (ii), after M updates the following two inequalities must hold:

- $\mathbf{w}^\top \mathbf{w}^* \geq M\gamma$ (1)
- $\mathbf{w}^\top \mathbf{w} \leq M$. (2)

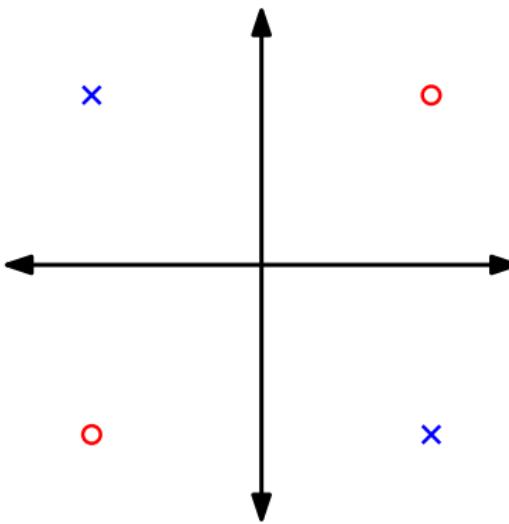
We can then complete the proof:

$$\begin{aligned} M\gamma &\leq \mathbf{w}^\top \mathbf{w}^* && \text{By (1)} \\ &= \|\mathbf{w}\| \cos(\theta) && \text{by definition of inner-product} \\ &\leq \|\mathbf{w}\| && \text{we must have } \cos(\theta) \leq 1. \\ &= \sqrt{\mathbf{w}^\top \mathbf{w}} && \text{by definition of } \|\mathbf{w}\| \\ &\leq \sqrt{M} && \text{By (2)} \end{aligned}$$

$$\Rightarrow M\gamma \leq \sqrt{M} \Rightarrow M \leq \frac{1}{\gamma^2} \quad \text{Hence, } M \text{ is bounded from above.}$$

A little bit history

- Initially, huge wave of excitement ("Digital brains") (See The New Yorker December 1958)
- Then, contributed to the first A.I. Winter. Famous example of a simple non-linearly separable data set, the XOR problem (Minsky 1969):



No linear separator classifies perfectly!

Neural networks via *features*.

To make a linear predictor nonlinear, we rely upon feature mapping ϕ :

$$\mathbf{x} \mapsto \mathbf{w}^\top \mathbf{x} \quad \text{becomes} \quad \mathbf{x} \mapsto \mathbf{w}^\top \phi(\mathbf{x}).$$

We are at the mercy of the quality of ϕ .

Why not *learn* ϕ ? e.g.,

$$\arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \ell \left(y^{(i)} \mathbf{w}^\top \mathbf{x}^{(i)} \right) \quad \text{becomes} \quad \arg \min_{\mathbf{w}, \phi \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell \left(y^{(i)} \mathbf{w}^\top \phi(\mathbf{x}^{(i)}) \right)$$

where \mathcal{F} is some class of functions (**why not every function?**).

Neural networks as iterated linear prediction (part 1).

Natural choice: build feature maps out of linear predictors!

$$\mathbf{x} \mapsto \mathbf{w}^\top \mathbf{x} \quad \text{becomes} \quad \mathbf{x} \mapsto \mathbf{v}^\top \phi(\mathbf{x}) \quad \text{where } \phi(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{b}$$

with $\mathbf{w} \in \mathbb{R}^d$, $\mathbf{v} \in \mathbb{R}^m$, $\mathbf{A} \in \mathbb{R}^{m \times d}$, $\mathbf{b} \in \mathbb{R}^m$.

There is something wrong with this!

Gained nothing! $\mathbf{v}^\top (\mathbf{A}\mathbf{x} + \mathbf{b}) = (\mathbf{A}^\top \mathbf{v})^\top \mathbf{x} + \mathbf{v}^\top \mathbf{b}$.

Fix: introduce **nonlinearity/transfer/activation** $\sigma : \mathbb{R}^m \rightarrow \mathbb{R}^m$:

$$\phi(\mathbf{x}) := \sigma(\mathbf{A}\mathbf{x} + \mathbf{b}).$$

Neural networks as iterated linear prediction (part 2).

We will predict with

$$\mathbf{x} \mapsto \mathbf{w}^\top \phi(\mathbf{x}) \quad \text{where } \phi(\mathbf{x}) = \sigma(\mathbf{A}\mathbf{x} + \mathbf{b}).$$

We will train with

$$\arg \min_{\mathbf{w} \in \mathbb{R}^m, \mathbf{A} \in \mathbb{R}^{m \times d}, \mathbf{b} \in \mathbb{R}^m} \frac{1}{n} \sum_{i=1}^n \ell \left(y^{(i)} \mathbf{w}^\top \sigma \left(\mathbf{A}\mathbf{x}^{(i)} + \mathbf{b} \right) \right).$$

(**Question:** which training procedure? Why does it work?)

Why stop there? We can also do

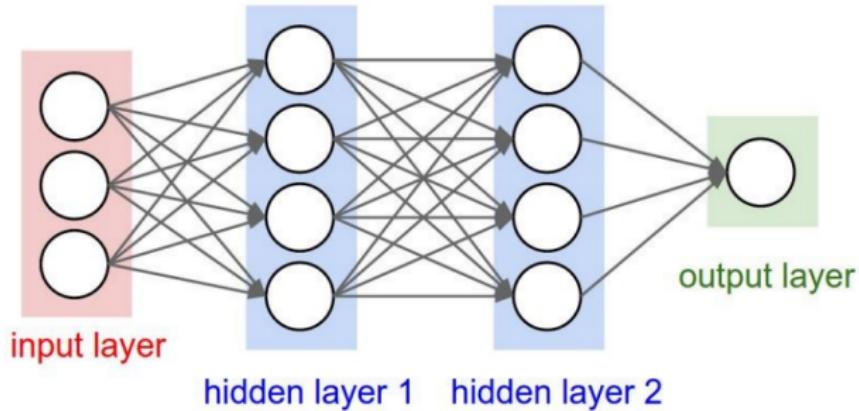
$$\mathbf{x} \mapsto \mathbf{w}^\top \sigma_1 (\mathbf{A}_1 \phi(\mathbf{x}) + \mathbf{b}_1) \quad \text{where } \phi(\mathbf{x}) = \sigma_2 (\mathbf{A}_2 \mathbf{x} + \mathbf{b}_2),$$

and iterate further. *This is a neural network.*

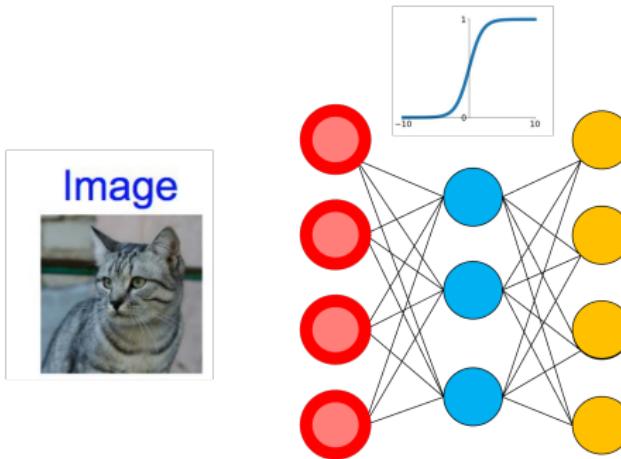
Multi-layer Perceptron

$$\mathbf{x} \mapsto \mathbf{w}^\top \sigma_1 (\mathbf{A}_1 \phi(\mathbf{x}) + \mathbf{b}_1)$$

where $\phi(\mathbf{x}) = \sigma_2 (\mathbf{A}_2 \mathbf{x} + \mathbf{b}_2)$,

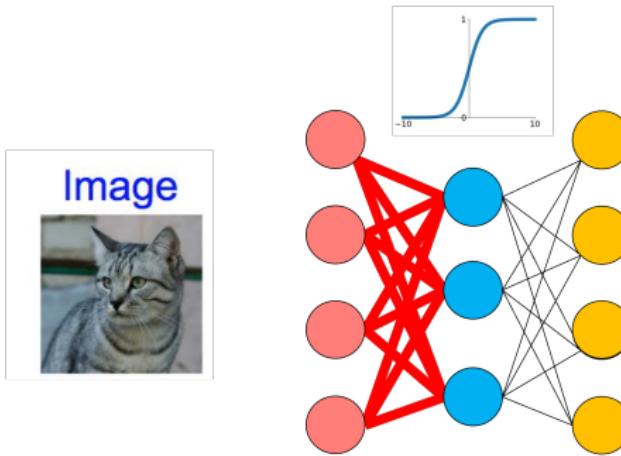


Forward pass



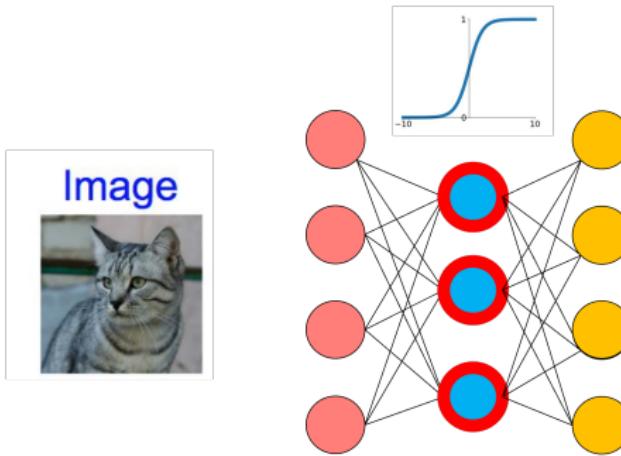
$$y = W_2 \sigma(W_1 x + b_1) + b_2$$

Forward pass



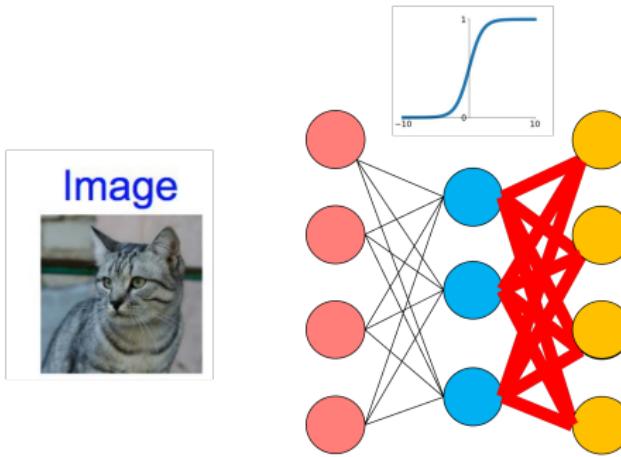
$$y = W_2 \sigma(W_1 x + b_1) + b_2$$

Forward pass



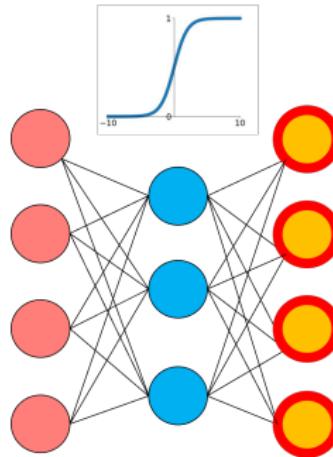
$$y = W_2 \sigma(W_1 x + b_1) + b_2$$

Forward pass



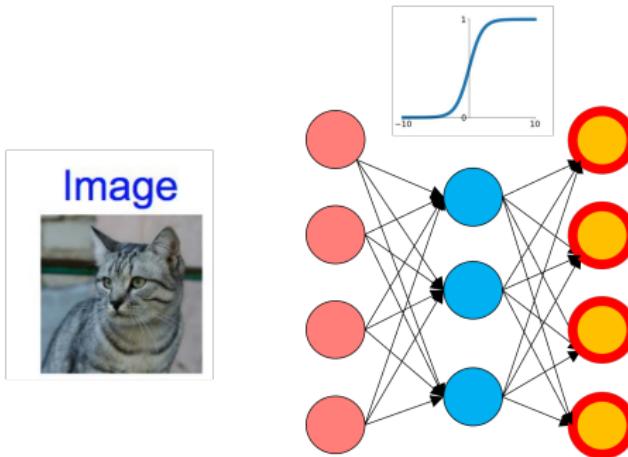
$$y = W_2 \sigma(W_1 x + b_1) + b_2$$

Forward pass



$$y = W_2 \sigma(W_1 x + b_1) + b_2$$

Forward pass



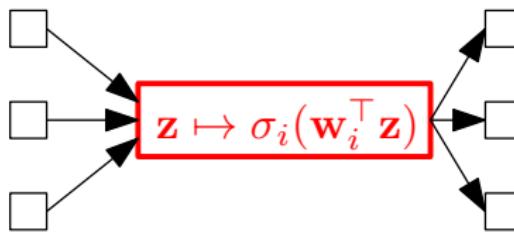
$$y = W_2 \sigma(W_1 x + b_1) + b_2$$

Classical formulation of neural networks as graphs.

(“Classical” because tensorflow “computation graphs” differ slightly.)

Node j in this graph:

- Collects a vector \mathbf{z} from its in-edges;
- Computes $\mathbf{z} \mapsto \sigma_j(\mathbf{w}_j^\top \mathbf{z} + b_j)$;
- Propagates this value along its out-edges.



Computation of whole network can be written this way.

Tensorflow computation graphs: everything needed to train is in the graph; e.g., parameters get nodes.

Neural networks as functions.

A linear predictor (**one layer network**) has the form

$$\mathbf{x} \mapsto \mathbf{w}^\top \mathbf{x}.$$

A **two layer network** has the form

$$\mathbf{x} \mapsto \mathbf{w}^\top \sigma_1(\mathbf{A}_1 \mathbf{x} + \mathbf{b}_1).$$

Iterating, a **multi-layer network** has the form

$$\mathbf{x} \mapsto \mathbf{w}^\top \sigma_1 \left(\mathbf{A}_1 \sigma_2 \left(\cdots \mathbf{A}_{L-2} \sigma_{L-1} \left(\mathbf{A}_{L-1} \mathbf{x} + \mathbf{b}_{L-1} \right) + \mathbf{b}_{L-2} \cdots \right) + \mathbf{b}_1 \right).$$

ERM now takes the form

$$\arg \min_{\mathbf{w}, \mathbf{A}_1, \dots, \mathbf{A}_{L-1}, \mathbf{b}_1, \dots, \mathbf{b}_{L-1}} \frac{1}{n} \sum_{i=1}^n \ell \left(y^{(i)} \mathbf{w}^\top \sigma_1 \left(\cdots \sigma_{L-1} (\mathbf{A}_{L-1} \mathbf{x}^{(i)} + \mathbf{b}_{L-1}) \cdots \right) \right).$$

Neural network (univariate) activations.

We mentioned that **nodes** compute

$$\mathbf{z} \mapsto \sigma(\mathbf{v}^\top \mathbf{z}),$$

where *activation/transfer/nonlinearity* $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is:

- ReLU (Rectified Linear Unit) $z \mapsto \max\{0, z\}$;
- Sigmoid $z \mapsto \frac{1}{1+\exp(-z)}$;
-

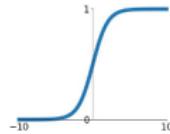
By $\mathbf{z} \mapsto \sigma(\mathbf{A}\mathbf{z} + \mathbf{b})$,
we meant “apply a univariate σ coordinate-wise”.

Soon we will see multivariate nonlinearities,
sometimes with output dimension \neq input dimension!

Nonlinearity.

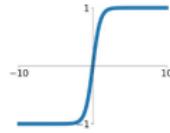
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



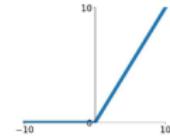
tanh

$$\tanh(x)$$



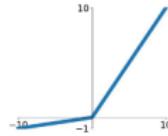
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

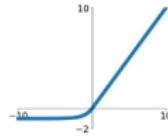


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Some modern usage.

Multiclass output.

Modern networks often end with **softmax** nonlinearity:

$$\mathbf{z} \mapsto \sum_{i=1}^k \frac{\exp(\mathbf{z}_i) \mathbf{e}_i}{\sum_{j=1}^k \exp(\mathbf{z}_j)}$$

(where \mathbf{e}_i is i^{th} standard basis vector.)

Output is now a probability vector!

Alternate notation: output vector $\mathbf{v}_i \propto \exp(\mathbf{z}_i)$.

Cross-entropy loss.

Given *one hot* $\mathbf{y} \in \{\mathbf{e}_1, \dots, \mathbf{e}_k\}$ and probability vector $\hat{\mathbf{y}} \in \mathbb{R}^k$,

$$\ell(y, \hat{y}) = - \sum_{i=1}^k \mathbf{y}_i \ln(\hat{\mathbf{y}}).$$

Combined with softmax $\hat{\mathbf{y}} \propto \exp(\mathbf{z})$:

$$-\sum_{i=1}^k \mathbf{y}_i \ln \left(\frac{\exp(\mathbf{z}_i)}{\sum_j \exp(\mathbf{z}_j)} \right) = -\sum_{i=1}^k \mathbf{y}_i \mathbf{z}_i + \ln \left(\sum_{i=1}^k \exp(\mathbf{z}_i) \right).$$

(For numerical stability, use $\ln \sum_i \exp v_i = c + \ln \sum_i \exp(v_i - c)$.)

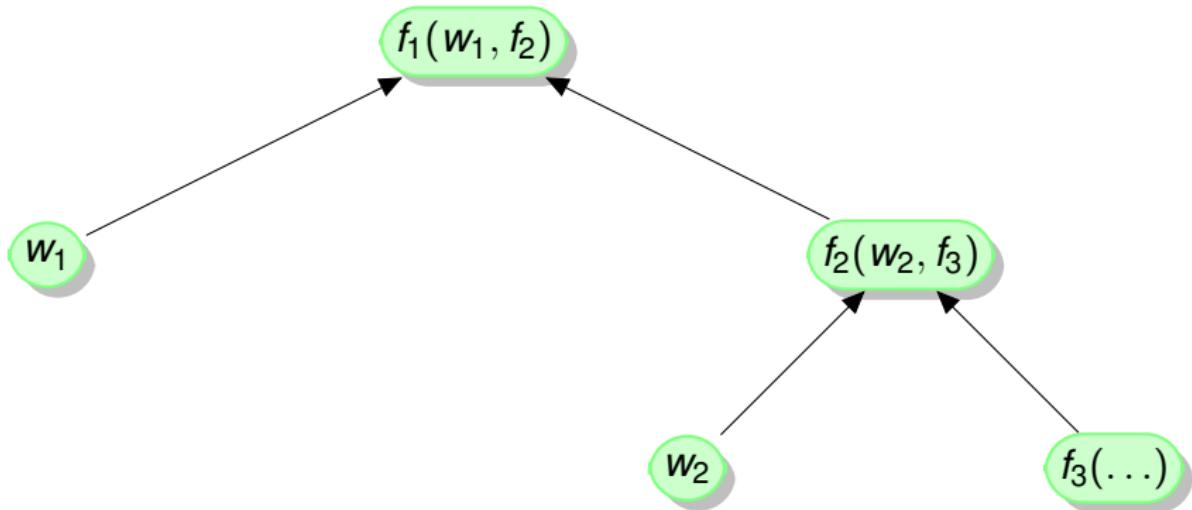
(In binary case $k = 2$: generalizes logistic loss.)

Question: since last expression is convex in \mathbf{z} ,
is the corresponding ERM problem convex?

Example:

$$F(\mathbf{w}, x, y) = f_1(w_1, f_2(w_2, f_3(\dots)))$$

Nodes are weights, data, and functions:



Internal representation used by deep net packages.

Example (PyTorch.py):

```
class Net(nn.Module):
    def __init__(self):
        super(Net, self).__init__()
        self.fc1 = nn.Linear(16, 120) % specify shape, init params
        self.fc2 = nn.Linear(120, 84) % specify shape, init params
        self.fc3 = nn.Linear(84, 10) % specify shape, init params

    def forward(self, x):
        x = F.relu(self.fc1(x)) % execute relu(w1*x + b1)
        x = F.relu(self.fc2(x))
        x = self.fc3(x)
        return x
```

What are the input dimensions?

Quiz:

- What are deep nets?
- How do deep nets relate to SVMs and logistic regression
- What components of deep nets do you know?

Important topics of this lecture

- Perceptron
- Multi-layer perceptron
- Components of deep nets

Up next:

- Backpropagation