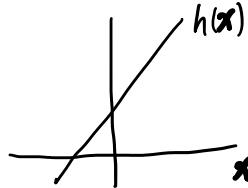


11-Nonlinear Models (Neural Networks)

Linear models

In previous topics, we mainly dealt with linear models

- Regression $h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$

$$h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$$


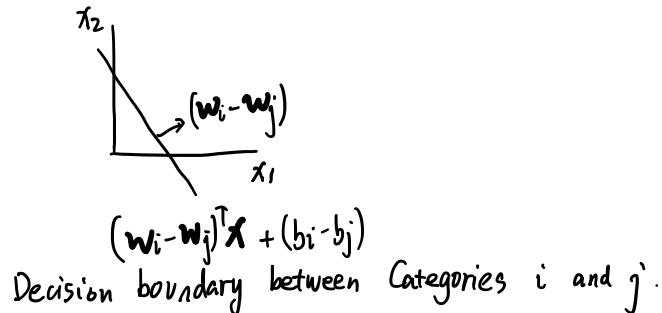
- Classification:

$$h(\mathbf{x}) = \operatorname{argmax}_i \mathbf{w}_i^\top \mathbf{x} + b_i$$

Category i is referred than Category j

$$\mathbf{w}_i^\top \mathbf{x} + b_i \geq \mathbf{w}_j^\top \mathbf{x} + b_j \quad \text{i.e., } (\mathbf{w}_i - \mathbf{w}_j)^\top \mathbf{x} + (b_i - b_j) \geq 0$$

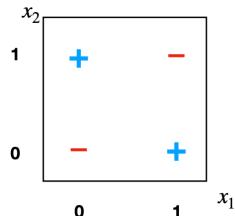
Binary classification with logistic regression is a special case of the above.



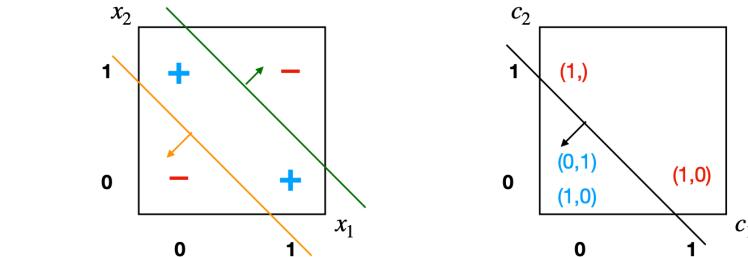
Nonlinear models

- Non-linear features $\phi(\mathbf{x})$
 - E.g., Gaussian discriminant analysis with the different covariance matrices, in which case we have quadratic features of \mathbf{x} .
- Non-linear kernel $k(\mathbf{x}_i, \mathbf{x}_j)$
 - A kernel is an inner-product of two data samples that are transformed in a certain vector space. The vector space could be very high-dimensional (e.g., with infinite dimensions). A linear classification in such a high-dimensional space could be non-linear in the original low dimensional space.
- Learnable non-linear mapping
 - We can probably stack a few layers of learnable non-linear functions (e.g., logistic functions) to learn the non-linear feature $\phi(\mathbf{x})$ or a non-linear kernel that is appropriate to the task at hand.

Motivation: XOR, a nonlinear classification problem

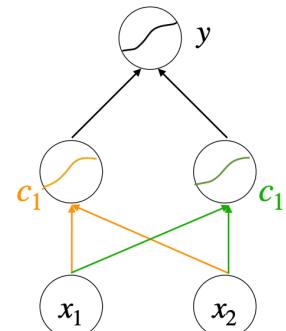


- The problem cannot be solved by logistic regression
- What if we stack multiple logistic regression classifiers?



Point:	(0,0)	(0,1)	(1,0)	(1,1)
Classifier 1:	1	0	0	0
Classifier 2:	0	0	0	1

- The XOR problem is solvable by three linear classifiers
 - One built upon the other two
 - But this is **programming** [HW]
 - Some machinery that allows you to specify certain things
 - Programming means you specify these things (usually heuristically by human intelligence) that are can be input to the machinery
 - The machinery accomplishes a certain task according to your input (program).
 - Programming is very tedious and only feasible simple tasks
 - We want to **learn** the weights.
- Can we learn the weights?
 - Yes, still by gradient descent.
- Can we compute the gradient?
 - Yes, it's still a differentiable function



$$y = \frac{1}{1 + e^{-(w_0 + w_1 c_1 + w_2 c_2)}}$$

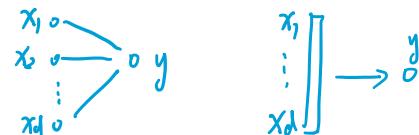
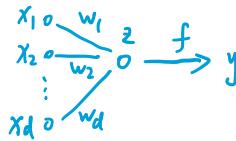
$$= \frac{1}{1 + e^{-(w_0 + w_1 \left\{ \frac{1}{1 + e^{-(w_{1,0} + w_{1,1}x_1 + w_{1,2}x_2)}} \right\} + w_2 \left\{ \frac{1}{1 + e^{-(w_{2,0} + w_{2,1}x_1 + w_{2,2}x_2)}} \right\}}}}$$

- Again, brute-force computation of gradient is very tedious.
- We need a systematic way of
 - Defining a deep architecture, and
 - Computing its gradient.

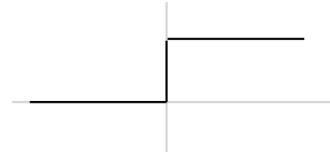
Artificial Neural Network

- A perceptron [Rosenblatt, 1958]

$$\begin{aligned} z &= \mathbf{w}^T \mathbf{x} + b \\ y &= f(z) \end{aligned}$$



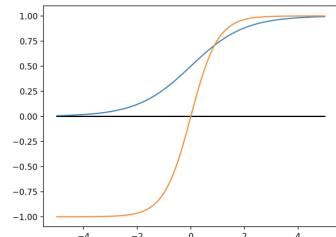
where f is a binary thresholding function



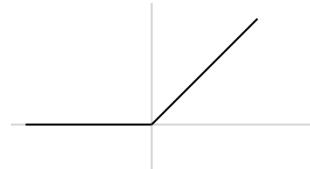
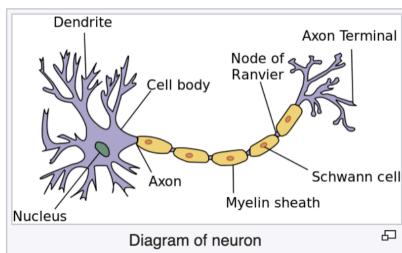
- A perceptron-like neuron, unit, or node

$$\begin{aligned} z &= \mathbf{w}^T \mathbf{x} + b \\ y &= f(z) \end{aligned}$$

f is an activation function, e.g., sigmoid, tanh, ReLU



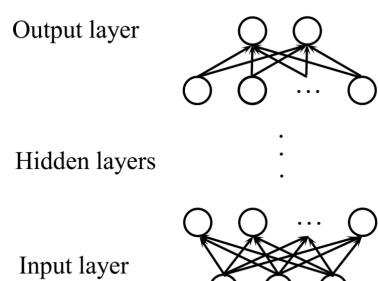
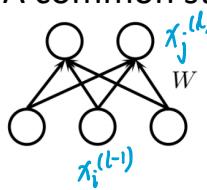
- Usually we use nonlinear activation
- Linear activation may be used for regression



Source: <https://en.wikipedia.org/wiki/Neuron>

- A multi-layer neural network, or a multi-layer perceptron

A common structure is layer-wise fully connected



For each node j at layer ℓ ,

$$z_j^{(\ell)} = \sum_i w_{ji}^{(\ell)} x_i^{(\ell-1)} + b_j^{(\ell)}$$

$$x^{(\ell)} = f(z^{(\ell)})$$

○ ○ ○

To simplify notations, we omit the layer L , but call the output of the current layer as y and the input of the current layer x , which is the output of the lower layer. In the simplified notation,

$z_j = \sum_i w_{ji} x_i + b_j$ $y_j = f(z_j)$ Node-wise notation	$z = Wx + b$ $y = f(z)$ Vector notation (common in papers)	$z = XW + \text{repmat}(b)$ $Y = f(z)$ Batch notation (Common in implementation)
	where $x \in \mathbb{R}^{N_x}$, $b, z, y \in \mathbb{R}^{N_y}$	where $X \in \mathbb{R}^{M \times N_x}$, $Z, Y \in \mathbb{R}^{M \times N_y}$
	$W \in \mathbb{R}^{N_y \times N_x}$	$W \in \mathbb{R}^{N_x \times N_y}$

Since we have multiple layers, we need a recursive algorithm that computes the activation of all nodes automatically.

Forward propagation (FP)

- Initialization $y^{(0)} = x$ known
- Recursion If $y^{(\ell-1)}$ is known, we can compute $y^{(\ell)}$
- Termination When $y^{(L)}$ is computed

Gradient of multi-layer neural networks

Main idea: if we can compute the gradient for one layer, we may use chain rule to compute the gradient for all layers.

Recursion on what?

Attempt #1 $\frac{\partial J}{\partial W^{(\ell)}}$: Bad idea. No direct connection between $\frac{\partial J}{\partial W^{(\ell)}}$ and $\frac{\partial J}{\partial W^{(\ell-1)}}$

Attempt #2: $\frac{\partial J}{\partial x^{(\ell)}}$ or $\frac{\partial J}{\partial z^{(\ell)}}$:

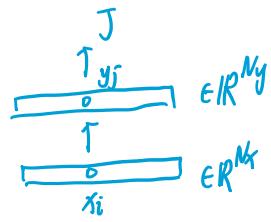
Good idea. $\frac{\partial J}{\partial \mathbf{x}^{(l-1)}}$ and $\frac{\partial J}{\partial \mathbf{W}^{(l)}}$ can be easily computed.

We consider a local layer

$$\text{FP: } \begin{cases} z_j = \sum_i w_{ji} x_i + b_j \\ y_j = f(z_j) \end{cases}$$

Gradient. Suppose $\frac{\partial J}{\partial y_j}$ is known

$$\begin{aligned} \frac{\partial J}{\partial x_i} &= \sum_j \frac{\partial J}{\partial z_j} \cdot \frac{\partial z_j}{\partial x_i} \\ &= \sum_j \underbrace{\frac{\partial J}{\partial y_j}}_{\substack{\text{recursive assumption} \\ (\text{f pointwise})}} \cdot \underbrace{\frac{\partial y_j}{\partial z_j}}_{w_{ji}} \cdot \underbrace{\frac{\partial z_j}{\partial x_i}}_{1} \end{aligned}$$



Note: $\frac{\partial J}{\partial z_j}$ is directly obtained at the output layer (MLE for GLIM)

$$\frac{\partial J}{\partial w_{ji}} = \frac{\partial J}{\partial z_j} \cdot x_i \quad \frac{\partial J}{\partial b_j} = \frac{\partial J}{\partial z_j}$$

$$\text{FP: } \begin{cases} z = Wx + b \\ y = f(z) \end{cases} \quad x \in \mathbb{R}^{N_x}, \quad y, z, b \in \mathbb{R}^{N_y}, \quad W \in \mathbb{R}^{N_y \times N_x}$$

$$\text{Assume } \frac{\partial J}{\partial y} = \begin{pmatrix} \frac{\partial J}{\partial y_1} \\ \vdots \\ \frac{\partial J}{\partial y_{N_y}} \end{pmatrix} \in \mathbb{R}^{N_y} \text{ known.}$$

$$\frac{\partial J}{\partial z} = \underbrace{\frac{\partial J}{\partial y}}_{\substack{\in \mathbb{R}^{N_y} \\ \in \mathbb{R}^{N_y}}} \cdot \underbrace{\left(\frac{\partial y_i}{\partial z_j} \right)}_{\substack{\in \mathbb{R}^{N_y} \\ \in \mathbb{R}^{N_y}}} \quad \cdot \cdot \cdot : \text{point-wise multiplication}$$

$$\frac{\partial J}{\partial x} = \underbrace{W^T}_{\substack{\in \mathbb{R}^{N_x} \\ \in \mathbb{R}^{N_y \times N_x}}} \underbrace{\frac{\partial J}{\partial z}}_{\substack{\in \mathbb{R}^{N_y} \\ \in \mathbb{R}^{N_y}}}$$

$$\frac{\partial J}{\partial W} = \underbrace{\frac{\partial J}{\partial z}}_{\substack{\in \mathbb{R}^{N_y} \\ \in \mathbb{R}^{N_y}}} \underbrace{x^T}_{\substack{\in \mathbb{R}^{N_x}}} \quad \frac{\partial J}{\partial b} = \frac{\partial J}{\partial z}$$

Batch representation = HW

Backpropagation (BP)

| o Initialization $\frac{\partial J}{\partial \mathbf{x}} \dots \frac{\partial J}{\partial \mathbf{z}}$

- Recursion If $\frac{\partial J}{\partial y^{(l)}}$ is known, then $\frac{\partial J}{\partial y^{(l-1)}}$ can be obtained
Consequently, we can compute $\frac{\partial J}{\partial w^{(l)}}$ and $\frac{\partial J}{\partial b^{(l)}}$
- Termination when $\frac{\partial J}{\partial w^{(l)}}$ and $\frac{\partial J}{\partial b^{(l)}}$ are computed

- A few more thoughts
 - Non-layerwise connection: Topological sort
 - Multiple losses: BP is a linear system
 - Tied weights: Total derivative is the summation

Auto-differentiation in general

- Input: <Layers, Edges, Losses>
- Algorithm:
 - Topological sort of all layers
 - Apply losses at respective layers
 - For layer L from last to first:
For each lower layer ℓ of L
 ℓ .gradient \doteq BP from L

Numerical gradient checking

- Definition of partial derivative

$$\frac{\partial}{\partial x_i} f(x_1, \dots, x_i, \dots, x_n) \stackrel{def}{=} \lim_{\delta \rightarrow 0} \frac{f(x_1, \dots, x_i + \delta, \dots, x_n) - f(x_1, \dots, x_i, \dots, x_n)}{\delta}$$

- Numerical gradient checking

$$\frac{\partial}{\partial x_i} f(x_1, \dots, x_i, \dots, x_n) \approx \frac{f(x_1, \dots, x_i + \delta, \dots, x_n) - f(x_1, \dots, x_i - \delta, \dots, x_n)}{2\delta}$$

