

Lecture 17

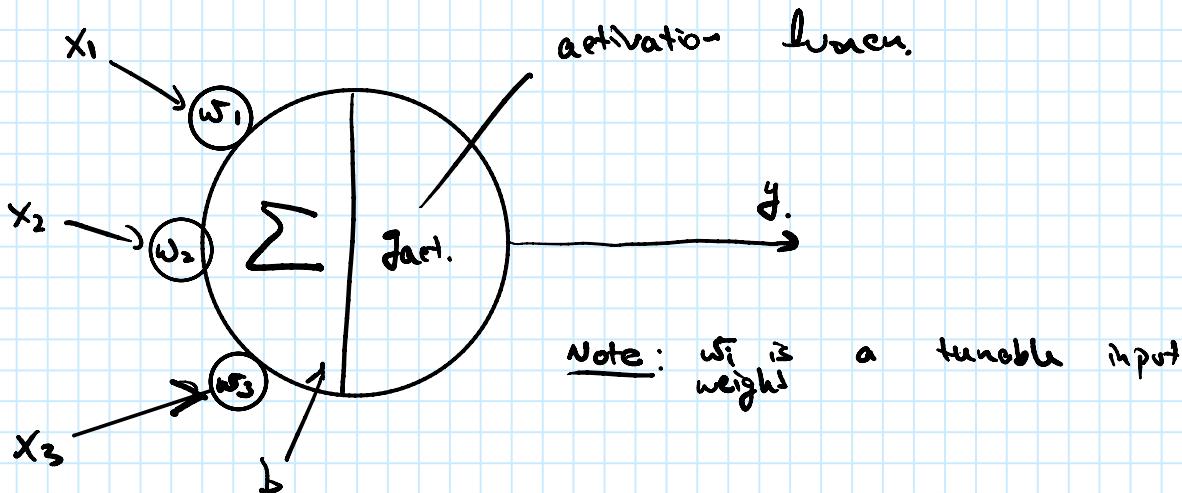
Wednesday, March 16, 2022 9:11 AM

Admin

- ① HW#3 due Friday.
② Exam Mar 30th.

Nonlinear Modeling - Part II

Last time: "Perceptron"



Key principle: - Given some dataset $\{x_1, \dots, x_n\}$ and a desired output y , the weights w_i and bias b can be chosen to create a "perceptron" that produces the same input-out behavior.

- The perceptron works well as a classifier for linearly separable data.
- The bias "b" allows the input threshold for the activation funen. to be adjusted to better fit the data.

• Activation functions for Artificial neurons:

- The activation can be written as:

$$g_{\text{act}}(z) \text{ where } z = \sum_{i=1}^n x_i w_i + b$$

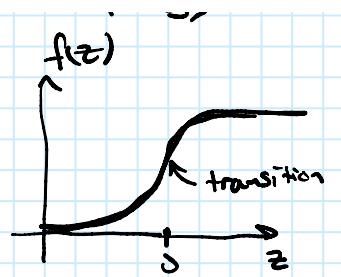
- $g_{\text{act}}(z)$ can take on several different forms, depending on the application.

$$f(z)$$

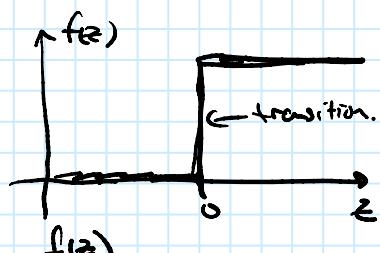
Depending on the application.

1) Sigmoid function: $f(z) = \frac{1}{1+e^{-z}}$

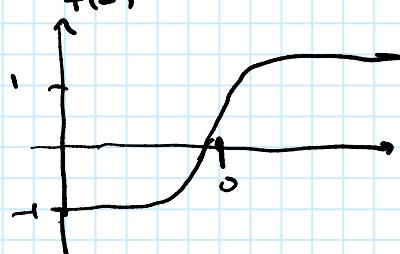
- generally outputs ≈ 0 , ≈ 1 but can have smaller outputs @ transition.



2) Threshold func: $f(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$

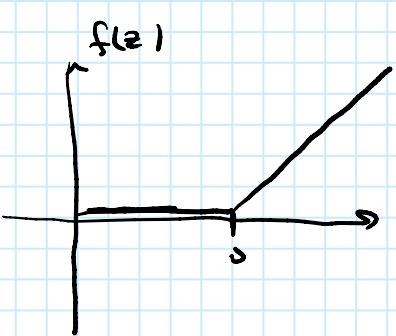


3) Hyperbolic tangent: $f(z) = \frac{1-e^{-2z}}{1+e^{-2z}}$

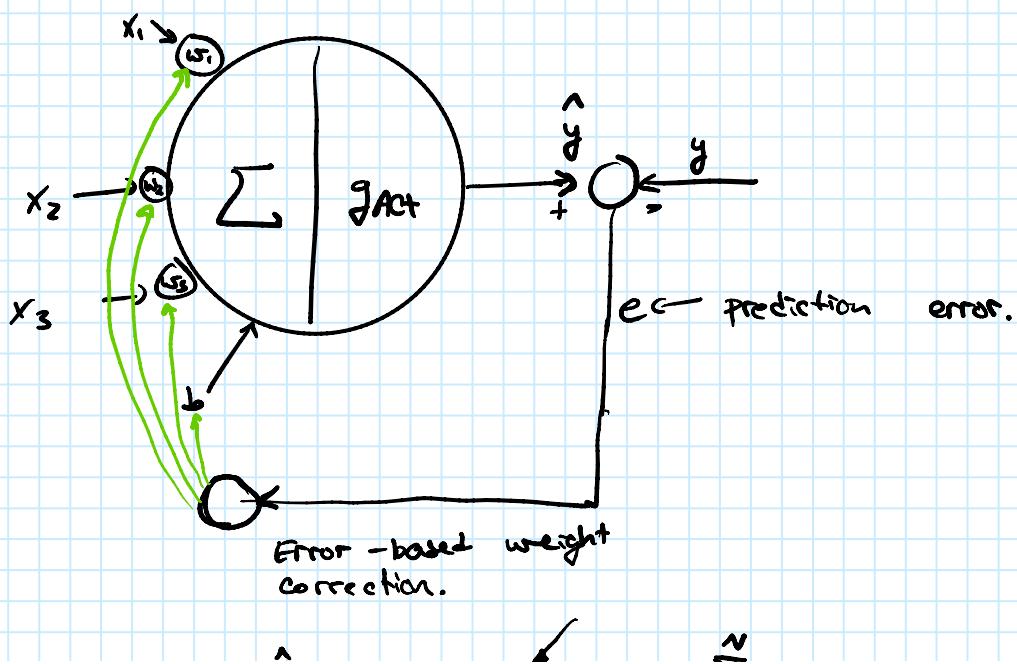


4) Rectified linear unit: $f(x) = \max(0, x)$

- Enables positive, proportional outputs.
- Also known as "ReLU"



- "Training" a perceptron by adjusting its input weights:



correction.

- Error term: $e = \hat{y} - y = g(z) - y$; $z = \sum_{i=1}^N w_i x_i + b$.

- One way to make corrections to use gradient descent

$$\Delta w_i = -\rho \underbrace{\text{grad}_{w_i} e^2}_{\substack{\text{"learning} \\ \text{rate"}}} = -\rho e \frac{\partial e}{\partial w_i}$$

↑
change in
 e w.r.t w_i

- The weight correction can be written as:

$$\Delta w_i = -2\rho \frac{\partial e}{\partial w_i} e \left(\frac{\partial e}{\partial w_i} \right) = x_i$$

↑
Assumes $g(z)$ is differentiable.

- This type of learning is referred to as "supervised learning".

- Used when you have a specific way you want the perceptron to respond to input data.

- For each data sample, j , provide (x_1^j, x_2^j, x_3^j) & y_j
⇒ Supervised learning rule: $\Delta w_i = (\text{Inputs } x_i) (\text{Error})$.

- There is also a type of training called "unsupervised learning".

- When you want the perceptron to create its own classification rules.

⇒ Unsupervised learning rule: $\Delta w_i \propto (\text{Input } x_i) (\text{output } y_i)$

⇒ see self organizing maps.

• procedures for training a perceptron:

- There are generally two ways to train:
one is a batch processing.

If we assume a linear activation function:

$$\hat{y} = g(z) = \sum_{i=1}^N w_i x_i + b$$

given sample data:

~ # samples.

given sample data:
 $(\hat{y}^j, x_1^j, \dots, x_n^j) | j=1, \dots, N$ # samples.

we want to find the weights to minimize

$$J_N = \frac{1}{N} \sum_{i=1}^N (\hat{y}^i - y^i)^2$$

\Rightarrow each weight can be computed by.

$$\Delta w_i = -\rho \text{grad}_{w_i} J_N = -\rho \frac{2}{N} \sum_j (\hat{y}^j - y^j) \frac{\partial \hat{y}^j}{\partial w_i}$$

\Rightarrow Requires that we store the gradient for all sample data $j=1, \dots, N$.

- The other way we can process is to update weights based on each new piece of data

For each new sample, K

$$\Delta w_i = \rho \delta^K x_i^K, \text{ where}$$

$$\delta^K = y^K - \underbrace{\sum_i w_i x_i^K}_\text{correct output} - b \quad \begin{matrix} \text{predicted output } \hat{y} \\ @ K^{\text{th}} \text{ presentation} \end{matrix} \quad \begin{matrix} \text{based on inputs} \\ \text{of data.} \end{matrix}$$

Learning procedure: Train perceptron by presenting all N training data in a segment, called an "epoch"



\rightarrow This process of presenting data multiple times is called "Recycling".

\rightarrow The learning procedure is called the Widrow-Hoff Algorithm.

- Convergence of Widrow-Hoff Algorithm.

- Convergence of Widrow-Hoff Algorithm.

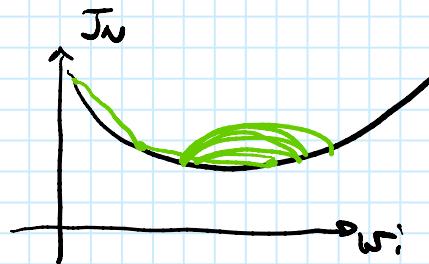
- If recycling is repeated infinite times, will the optimal weight vector be found?

$$\Rightarrow \arg\min J_W(w_1, w_2, \dots, w_n)$$

- Answer: If a constant learning rate $\rho > 0$ is used, the algorithm only converges if $\min J_W = 0$

If $J_W \neq 0$, Δw_i will also be nonzero for each epoch.

$\Rightarrow J_W$ will not settle.



- \Rightarrow To prevent this, we can use a variable learning rate:

$$\rho_k = \frac{\text{constant}}{k}$$



\Rightarrow This damps oscillations as k increases, improving convergence.