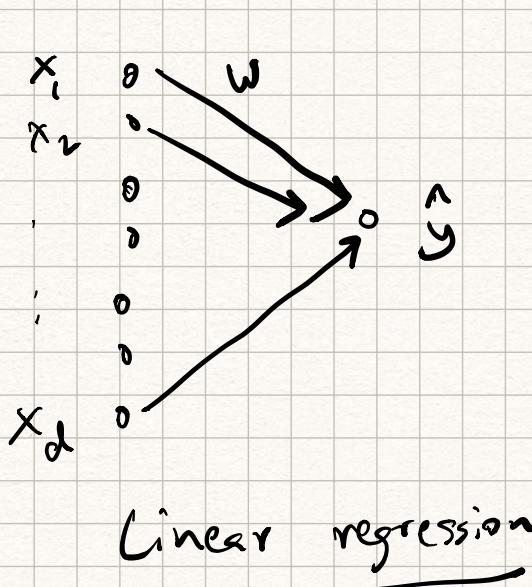


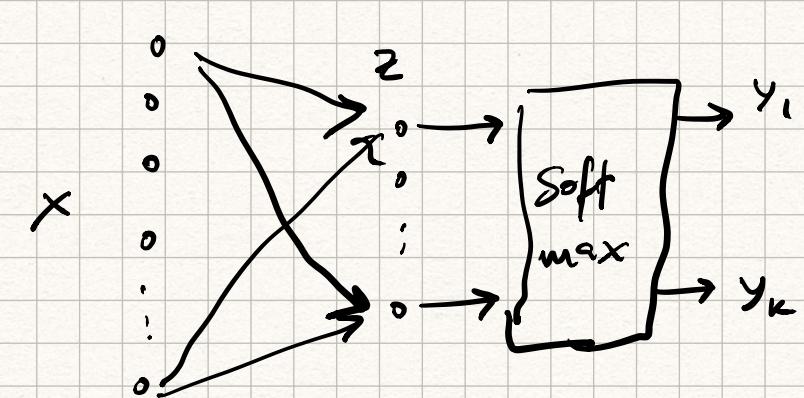
ECE-GV 6143  
Lecture 9

- Recap: linear models
- Neural nets
  - Architecture
  - Training via backpropagation.

- Reminder : Midterm exam due at 8 pm ET.



Linear regression



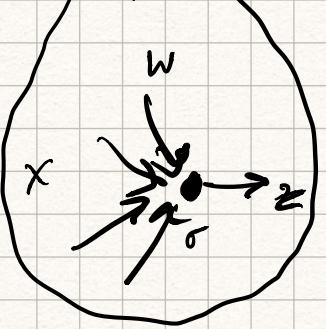
Logistic regression

- Feed forward
- Directed acyclic graphs

Neural networks

Primitive:  $z = \sigma(\langle w, x \rangle + b)$

"Neuron".

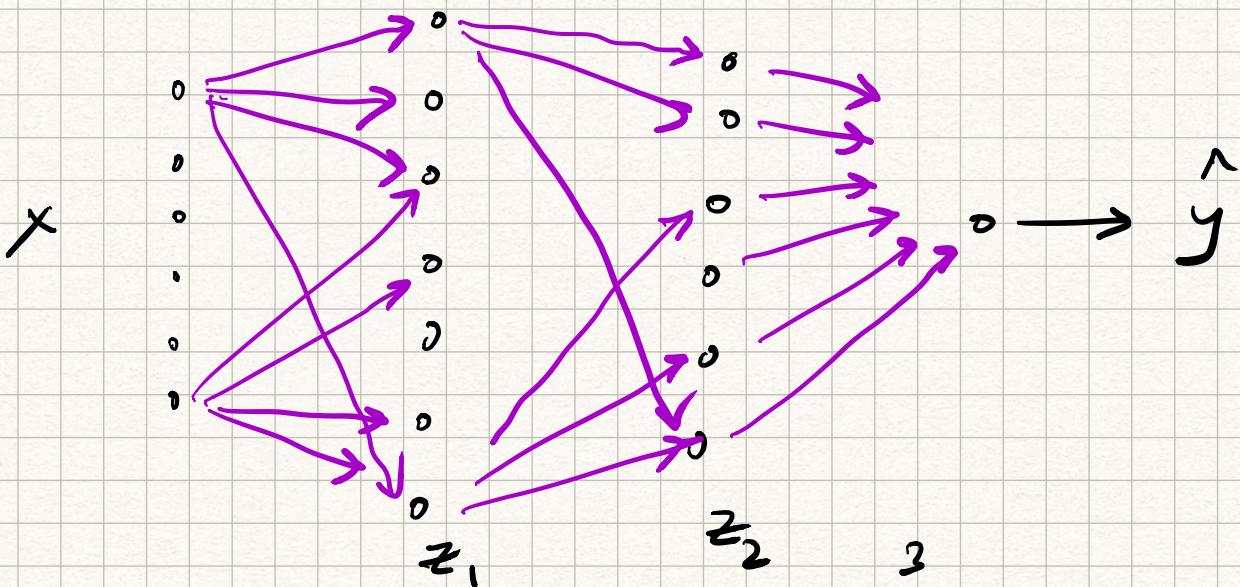
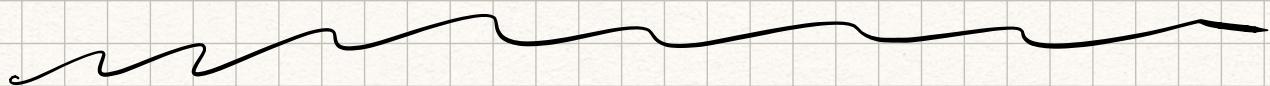


$\sigma$ : "Activation function".

Examples:  $\sigma(z) = z$   $\rightarrow$  Linear regression

$$\sigma(z) = \frac{1}{1+e^{-z}} \rightarrow \text{Logistic regression.}$$

$$\sigma(z) = \text{sign}(z) \rightarrow \text{Perception}$$



$$z_1 = \sigma^{(1)}(w^{(1)}x + b^{(1)})$$

$$z_2 = \sigma^{(2)}(w^{(2)}z_1 + b^{(2)})$$

$$\hat{y} = \sigma^{(3)}(w^{(3)}z_2 + b^{(3)}).$$

Functional form

"3-layer network" / "2-hidden-layer network"

Why?

- Historical / biological
- Mathematical.

Universal approximation theorem (Cybenko 1989)

:= Every continuous function  
can be approximated by a 1-hidden  
layer network of finite size (for  
most activation functions)

- Drawback

- Existence theorem only
- Practical issues

Representation □ What network size & shape?  
Training □ How to choose the network weights?

Generalization □ Does the learned generalize  
to unseen inputs?

Neural network architecture.

Design criteria : 1) Weights

## 2) Activation function.

### Activation functions

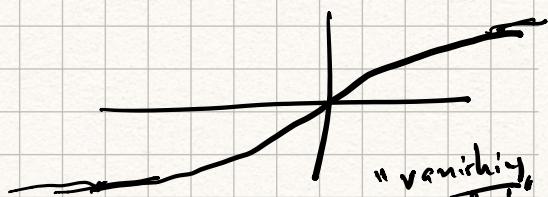
- $r(z) = z$

[Not too useful for multi-layer networks]

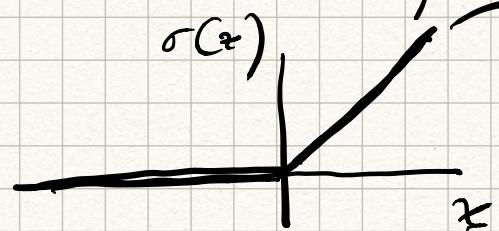
$$y = \frac{w^{(3)} w^{(2)} w^{(1)} x}{w_x}$$

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

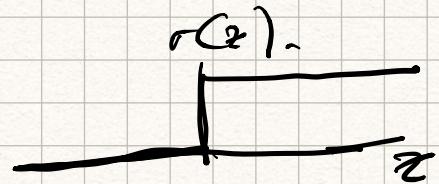
- $\sigma(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$  (tanh)



- $\sigma(z) = \text{ReLU}(z) = \max(z, 0)$



- $\sigma(z) = \text{sgn}(z)$



- $\delta(z) = HT(z)$

### Weights.

- Dense layers  $\rightarrow$  weights are arbitrary.
- Convolution layers  $\rightarrow$  weights implement convolution.
- Pooling layers  $\rightarrow$  Downsampling size of output.

- Batch normalization  $\rightarrow$  Rescaling.
- Recurrent layers  $\rightarrow$  Feed back.
- Attention layers  $\rightarrow$  NLP (etc.)

Q: How do I mix & match?

A: No correct answer.

Thumb rule: Just use a good existing architecture..



## Training

Given architecture, how do I train a network?

- \* Define loss function (+ optional regularization).
- \* Use some variant of gradient descent ( $GD / SGD / \dots$ ) -

$\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$

$$-L(W) = \frac{1}{n} \sum_{i=1}^n l(y_i, f_W(x_i)) + R(W)$$

$$-W^{t+1} \leftarrow W^t - \alpha^t \nabla L(w)$$

toy example

$$x \xrightarrow{\sigma} y$$

$$z = wx + b$$

$$f(z) = \sigma(z)$$

$$L(w, b) = \frac{1}{2} (y - f(z))^2 + \lambda w^2$$

$$= \frac{1}{2} (y - \sigma(wx + b))^2 + \lambda w^2$$

$$\nabla L(w, b) = \begin{bmatrix} \frac{\partial L}{\partial w} \\ \frac{\partial L}{\partial b} \end{bmatrix}$$

$$\frac{\partial L}{\partial w} = (\underbrace{\sigma(wx + b) - y}_{+ 2\lambda w}) \cdot \underbrace{\sigma'(wx + b)}_{\cdot x}$$

$$\frac{\partial L}{\partial b} = (\underbrace{\sigma(wx + b) - y}_{+ 2\lambda w}) \underbrace{\sigma'(wx + b)}_{\cdot x}$$

Backpropagation / backprop-

- exploits the structure of the network

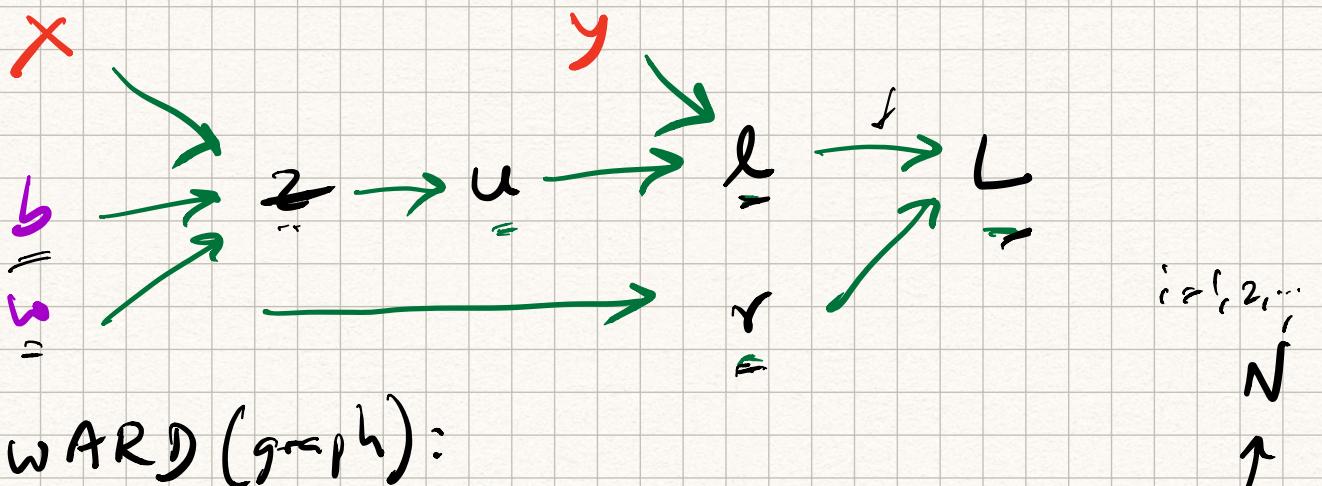
$$z = wx + b$$

$$u = \sigma(z)$$

$$l = \frac{1}{2} (u - y)^2$$

$$r = w^2$$

$$L = l + \lambda \sigma.$$



FORWARD (graph):

For node  $i = 1, 2, \dots, N$  :

compute  $v_i$  as a function  
of parents of  $i$ .

Number  
of nodes

BACKWARD (graph):

For node  $i = N-1, \dots, 1$  :

Compute

$$\frac{\partial L}{\partial v_i} = \sum_{j \in \text{children}(i)} \frac{\partial L}{\partial v_j} \cdot \frac{\partial v_j}{\partial v_i}$$

$$\frac{\partial L}{\partial v} := \frac{\partial L}{\partial v}$$

$$\frac{\partial L}{\partial e} := 1$$

$$\frac{\partial_e L}{\partial e} := 1$$

$$\begin{aligned}
 \partial_r L &= \lambda \\
 \partial_u L &= \partial_e L \cdot \partial_u^e l \\
 &\quad = 1 \cdot (u - y) \\
 \underline{\partial_x L} &= \partial_u L \cdot \underline{\partial_z^u} \\
 &\quad = \partial_u L \cdot \sigma'(x) \\
 \left\{ \begin{array}{l} \partial_w L = \underline{\partial_x L} \cdot \partial_w^x + \partial_r L \cdot \partial_w^r \\ \partial_b L = \underline{\partial_x L} \cdot x + \lambda \cdot 2w \end{array} \right. \\
 &\quad = \underline{\partial_x L} \cdot
 \end{aligned}$$

- Computationally efficient
  - Modular
  - Scalable.
-