

# Feed-forward Neural Networks

$$\underline{f} = g^{(3)} \left( \underline{W}^{(3)} \underline{h}^{(2)} + \underline{b}^{(3)} \right)$$

$$\uparrow$$
$$\underline{h}^{(2)} = g^{(2)} \left( \underline{W}^{(2)} \underline{h}^{(1)} + \underline{b}^{(2)} \right)$$

$$\uparrow$$
$$\underline{h}^{(1)} = g^{(1)} \left( \underline{W}^{(1)} \underline{x} + \underline{b}^{(1)} \right)$$

$\uparrow$   
 $\underline{x}$ , inputs

---

$$h_k^{(2)} = g^{(2)} \left( \sum_l W_{kl}^{(2)} h_l^{(1)} + b_k^{(2)} \right)$$

Homework: Try plotting for random weights.

Change: scale  $W$ , # layers,  $g$ , without  $\underline{b}$

Fitting params  $\theta$   $\rightarrow$  Vector containing  $\rightarrow$  # layer

Stochastic Gradient Descent  $\{W^{(l)}, b^{(l)}\}_{l=1}^L$   
on cost  $c$  and any other params.

$$\underline{\theta} \leftarrow \underline{\theta} - \eta \nabla_{\theta} c(\underline{\theta})$$

$\eta$  - learning rate

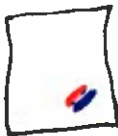
Tomorrow!  
how to compute  $\nabla_{\theta} c(\theta)$

Cost  $c$  describes task

- Match  $f$  to one-hot class label
- Match  $f$  to strength of concrete
- Match  $f$  to pixels of next image

...

Side-effect: learn "hidden" representations  $h^{(l)}$ , and  $e$

$W_{k,:}^{(1)}$  as image:  this "filter" looks for a type of edge.

$h^{(1)}$  indicates which edges present.

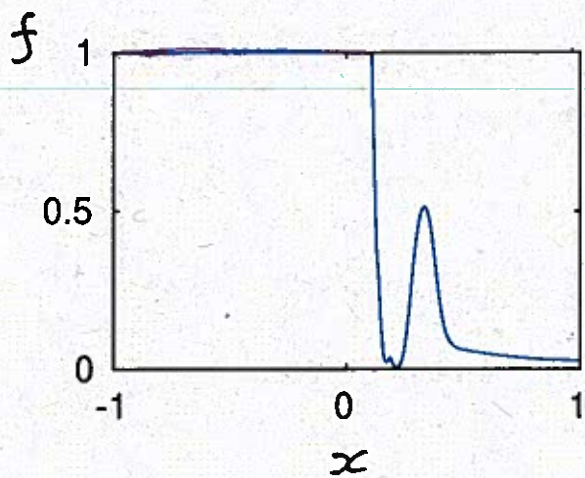
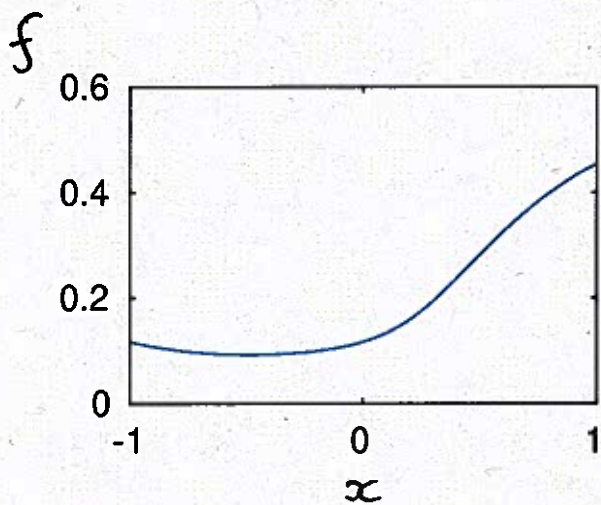
Embeddings:

If  $x$  one-hot,  $W_{:,d}^{(1)}$  is vector for choice  $d$

[Often set  $h^{(1)}$  directly to  $W_{:,d}^{(1)}$   
no point doing  $W_x^{(1)}$  and applying function]

```
gg = @(a) 1./(1 + exp(-a));  
X = (-1:0.01:1)'; % Nx1  
H1 = gg(X * randn(1, 100)); % Nx100  
H2 = gg(H1 * randn(100, 50)); % Nx50  
F = gg(H2 * randn(50, 1)); % Nx1  
plot(X, F);
```

```
gg = @(a) 1./(1 + exp(-a));  
X = (-1:0.01:1)'; % Nx1  
H1 = gg(X * randn(1, 100) * 10); % Nx100  
H2 = gg(H1 * randn(100, 50) * 10); % Nx50  
F = gg(H2 * randn(50, 1) * 10); % Nx1  
plot(X, F);
```



# Specialized Architectures

ConvNets: learned filters  move over whole image  
don't reload in every location  
⇒ share parameters  
...

## Simple pooling (example)

"Bag of words" text classification

$$\underline{e}^{(t)} = W \underline{x}^{(t)} = \text{"embedding vector for } t^{\text{th}} \text{ word"}$$

↳ onehot

$$\underline{h} = \left(\frac{1}{T}\right) \sum_{t=1}^T \underline{e}^{(t)} \quad \text{"document vector"}$$

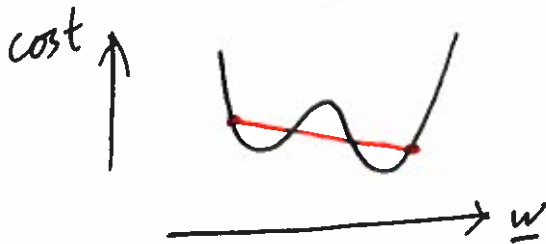
↳ ?

$$\underline{f} = \text{Softmax}(V \underline{h})$$

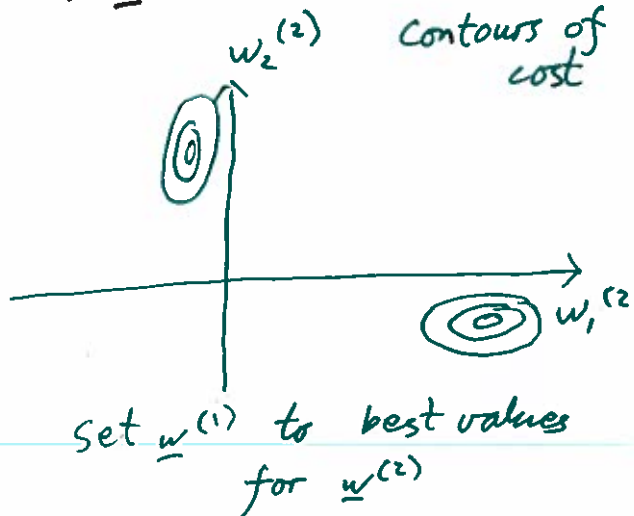
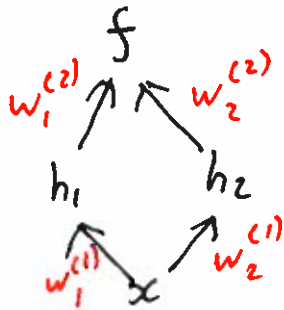
↳ prob. over classes

↳ embedding dim × |  
# classes × embedding dim

# NN don't have a convex cost



No unique optimum  
 $\Rightarrow$  not convex



- These local optima, don't matter

$\rightarrow$  because the functions (predictions) are same.

- Not all optima are equivalent.

## Initialization

Must not set all  $w^{(l)}$  to zero

→ all hidden extract same feature

→ SGD applies same update to all  
params

→ all hidden stay same...

Set randomly  $w_{ij}^{(l)} \sim N(0, 1)$  ?

[Cf Wof note]

- activation  $w_{\underline{x}}^{(l)}$  sum of  $D$  random  
values for each  
hidden

- typically each term in sum  
 $\sim \pm 1$  (if data  
standardized)

-  $\underline{w}^T \underline{x}$  is  $\sim \pm \sqrt{D}$

- instead  $w \sim N(0, (\frac{1}{\sqrt{D}})^2)$

MLP course notes have more ideas.