

Feed-forward Neural Networks

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

↑

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

↑

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

↑

\underline{x} , inputs

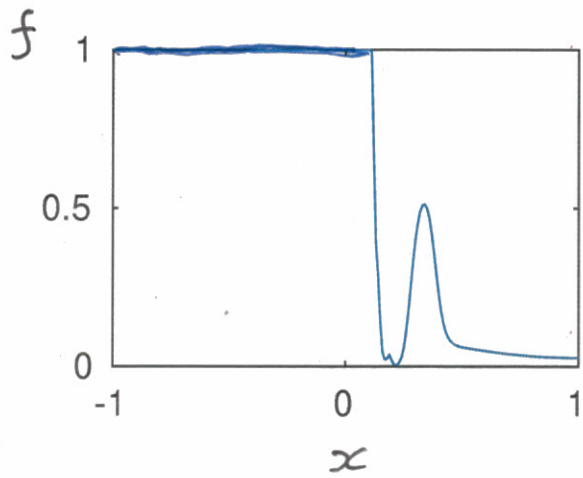
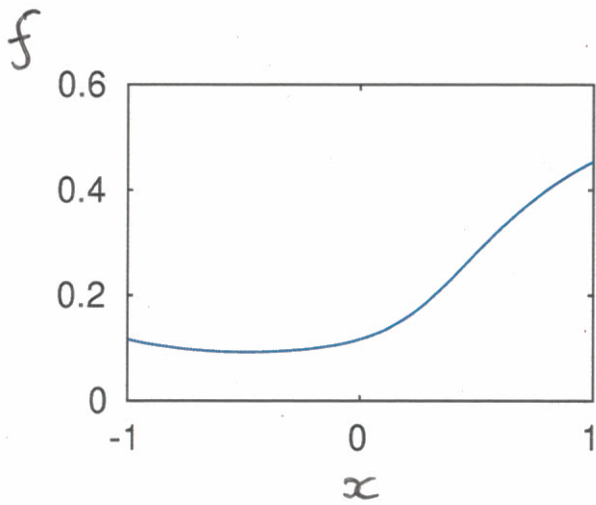
$$h_k^{(2)} = g^{(2)} \left(\sum_l W_{kl}^{(2)} h_l^{(1)} + b_k^{(2)} \right)$$

Homework: Try plotting for random weights.

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

```
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);
```



Initialization

Don't set all weights ^{to zero.} or the same value

Name $W_{ij}^{(l)} \sim N(0, 1)$

Use "randn"

[Revise background notes expectations]

Sum of k random values

They are typically ± 1 each

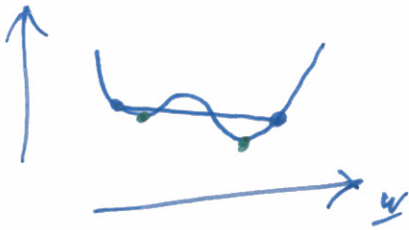
Typically sum $\sim \pm \sqrt{k}$

Example initialization $w \sim N(0, (\frac{1}{\sqrt{k}})^2)$

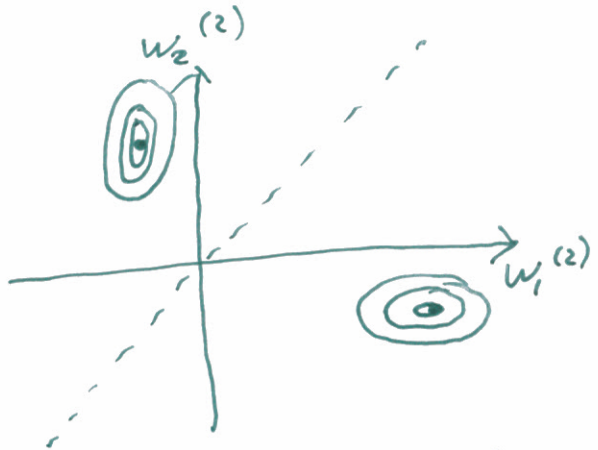
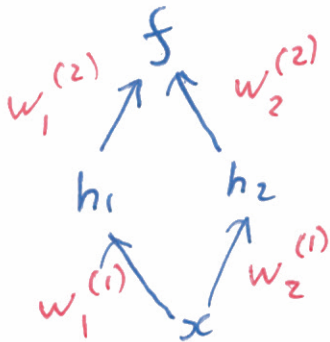
MLP: Has more sophisticated ideas.

NN don't have a convex cost

cost



No unique optimum
⇒ not convex.



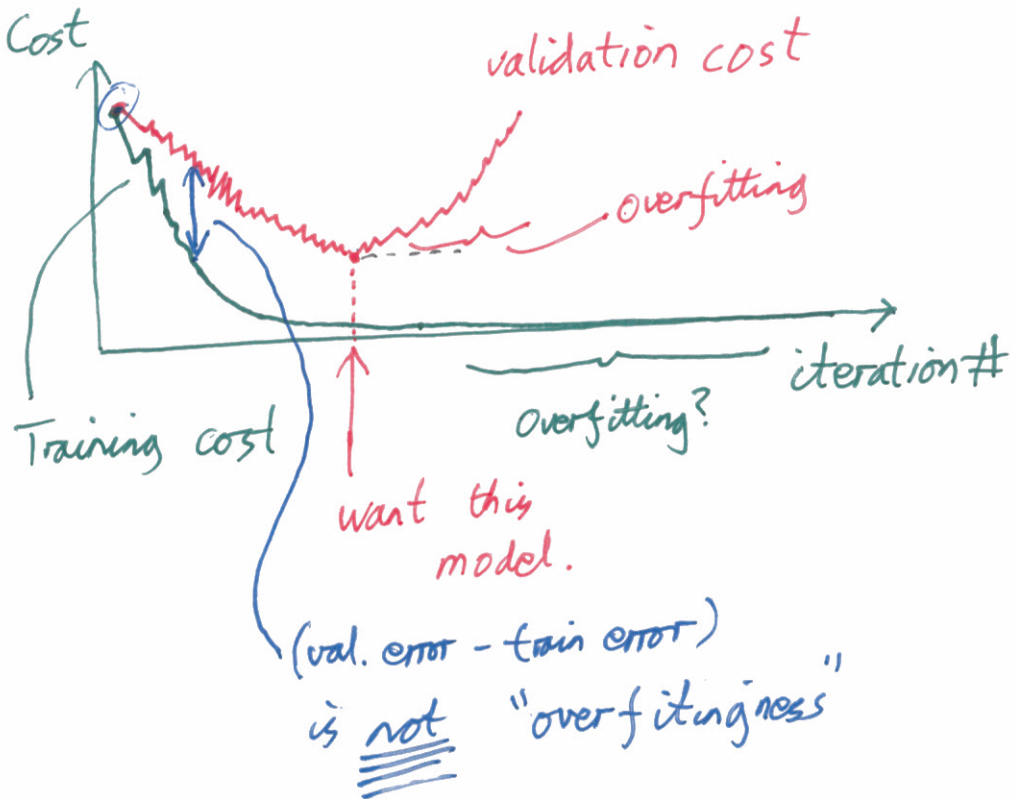
Set $\underline{w}^{(1)}$ to best values
for $\underline{w}^{(2)}$

- These local optima, don't matter
→ because the functions (predictions)
are the same.
- Not all optima are equivalent

Regularization - Early-stopping

Could do L2 regularization

Set λ ... cross-validate



Every R updates:

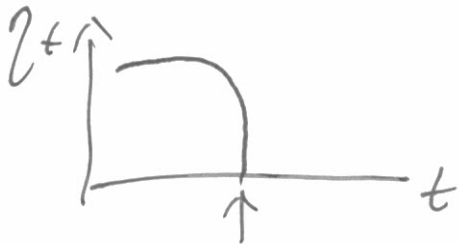
If val. cost is the smallest I've seen:

Store: weights & val cost.

If val. cost hasn't improved in 20 ~~updates~~

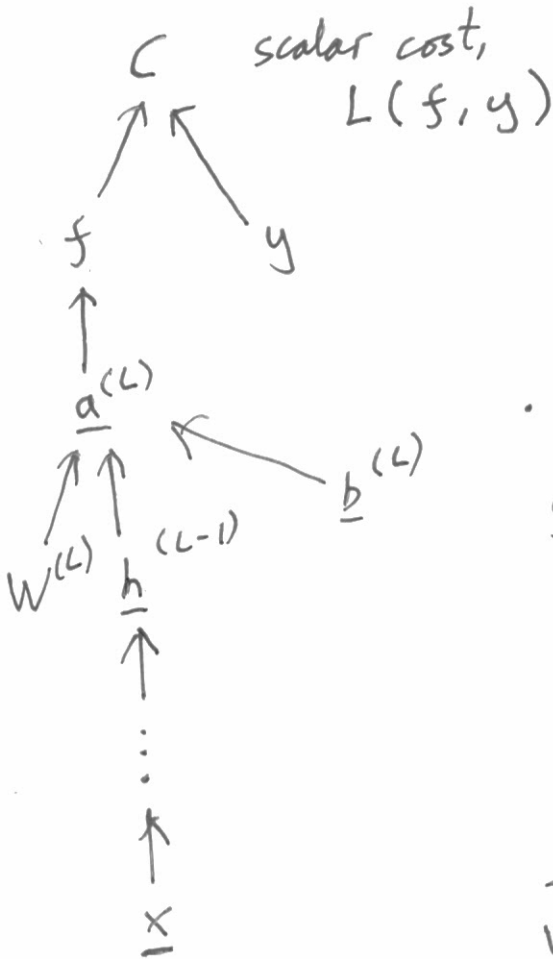
Stop. Return weights ^{at} _{val} ^{looks} _{at} ^{val} _{val} last stored.

Learning rates η_t → or reduce η



Getting gradients - Reverse-mode differentiation

"Back propagation"



Graph
DAG

Strategy:

- For every intermediate Θ

$$\text{get } \bar{\Theta} = \frac{\partial C}{\partial \Theta}$$

$$\bar{f} = \frac{\partial C}{\partial f}$$

$$\bar{a}_i^{(L)} = \frac{\partial C}{\partial a_i^{(L)}}$$

$$\bar{W}_{ij}^{(L)} = \frac{\partial C}{\partial W_{ij}^{(L)}}$$