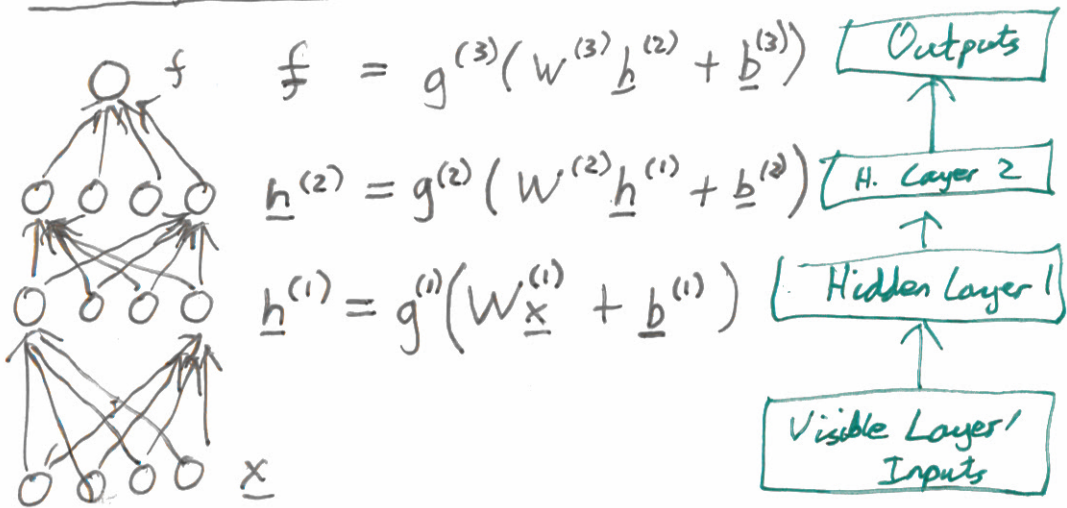


Feed-forward Neural Networks



- When f is a scalar, $W^{(3)} \underline{h}^{(2)} = \underline{w}^{(3)T} \underline{h}^{(2)}$

- Other architectures possible:

"skip connections"

parameterize the g 's non-linearities

- Special layers for images/audio

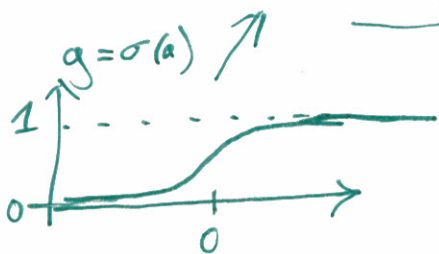
Conv Nets ... and others.

Initialization

Don't set all weights the same.

Naive code $W_{ij}^{(l)} \sim N(0, 1)$ "use randn()"

$$h_k^{(l)} = g \left(\underbrace{W_{k:}^{(l)} \underline{h}^{(l-1)}}_u + b_k^{(l)} \right)$$



Typically how big is this sum over $K^{(l-1)}$

Background on expectations:

Typically sum $\sim \pm \sqrt{K^{(l-1)}}$

Example Initialization $w \sim N(0, (1/\sqrt{K})^2)$

MLP: More sophisticated suggestions.

Showed MLPR 2017 U2 ③
Again

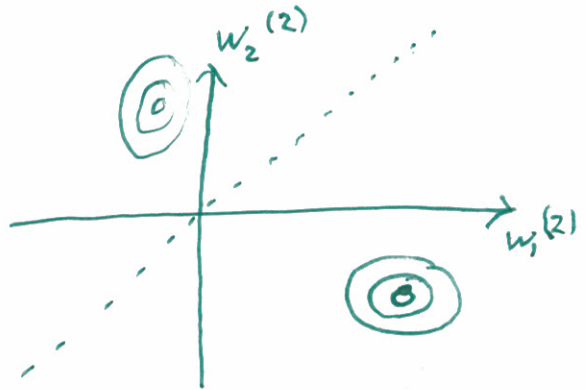
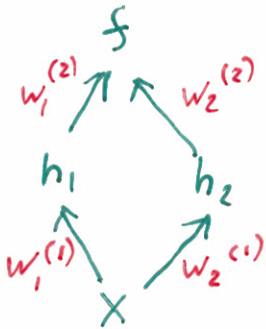
MLPR 2017 U3 ①

NN don't even have a convex cost



Convex

⇒ Unique optimum



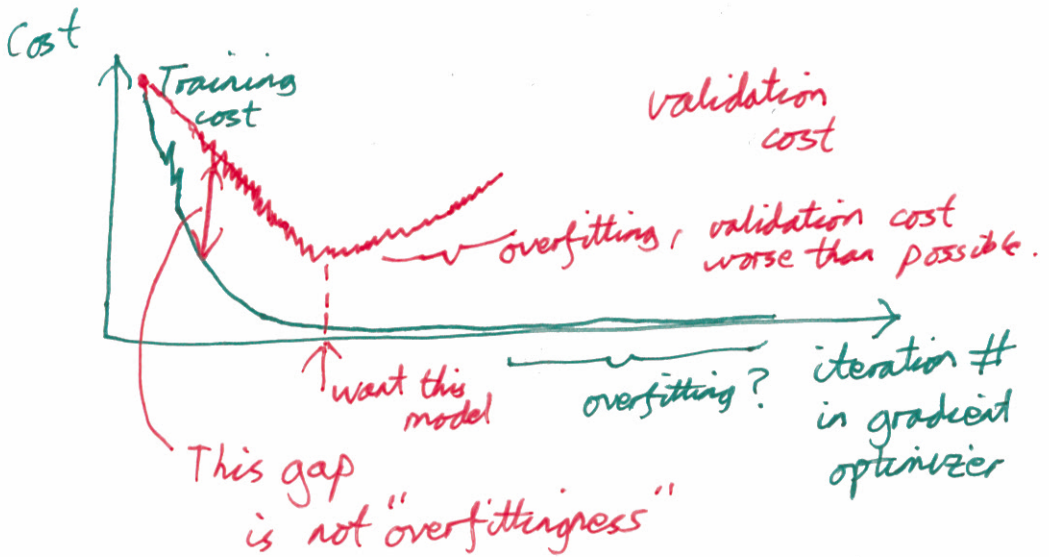
Set $\underline{w}^{(1)}$ to their best values.

- These local optima don't matter
 - because the function (predictions) are the same.
- Not all optima are equivalent.
- Use heuristics to get good fits.
- Random restarts. - not if network large/expensive

Regularization

Could do L2 regularization

Set $\lambda \dots$ cross-validate.



Every k updates:

If val. cost the smallest I've seen:

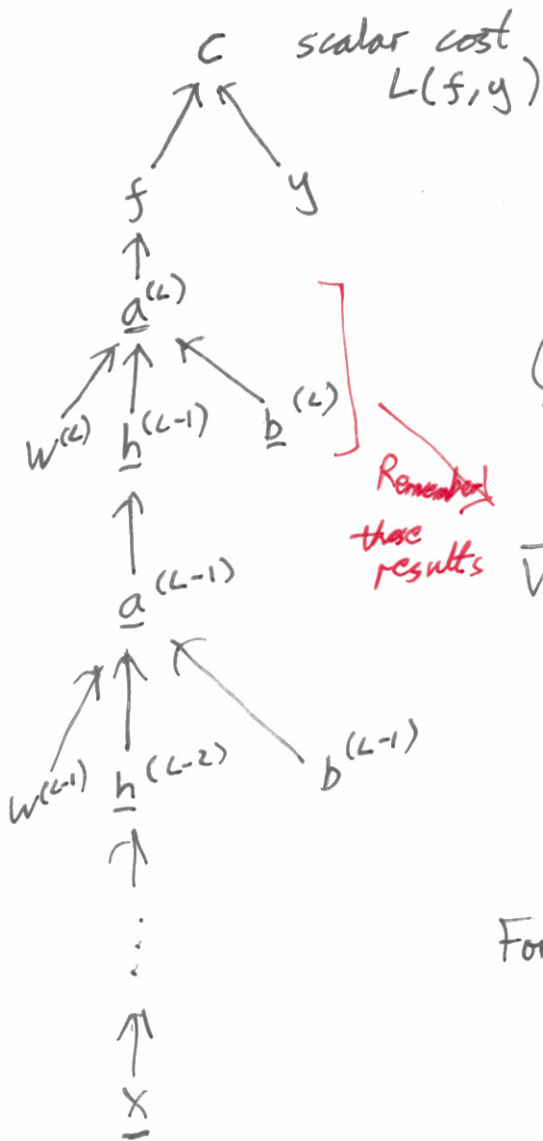
Store the weights & val. cost

If val. cost hasn't improved in 20 updates

Stop. Return the weights we stored from best val. score.

Getting gradients - Reverse-mode differentiation

Backpropagation



Remember
these
results

Strategy:

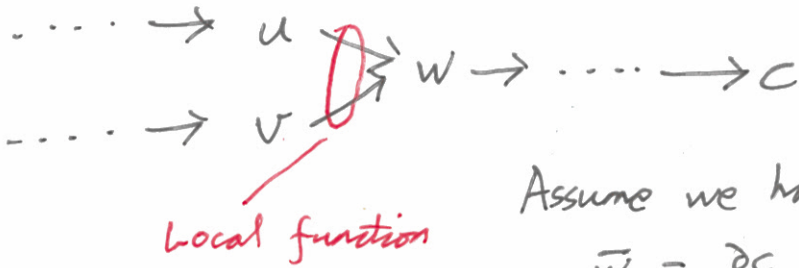
$$\bar{f} = \frac{\partial c}{\partial f}$$
$$\bar{a}_i^{(L)} = \frac{\partial c}{\partial a_i^{(L)}}$$
$$\bar{w}_{ij}^{(L)} = \frac{\partial c}{\partial w_{ij}^{(L)}}$$

Diagram illustrating the backward pass (backpropagation) of gradients. The gradient \bar{f} is calculated as $\frac{\partial c}{\partial f}$. The gradient $\bar{a}_i^{(L)}$ is calculated as $\frac{\partial c}{\partial a_i^{(L)}}$. The gradient $\bar{w}_{ij}^{(L)}$ is calculated as $\frac{\partial c}{\partial w_{ij}^{(L)}}$. The diagram shows the flow of gradients from the output f back through the layers L , $L-1$, and $L-2$, and from the target y to the output f .

For every intermediate θ

$$\text{get } \bar{\theta} = \frac{\partial c}{\partial \theta}$$

Derivative Propagation



Assume we have

$$\bar{w} = \frac{\partial c}{\partial w}$$

Want:

$$\bar{u} = \frac{\partial c}{\partial u}, \bar{v} = \frac{\partial c}{\partial v}$$

Chain rule:

$$\begin{aligned} \frac{\partial c}{\partial u} &= \frac{\partial c}{\partial w} \frac{\partial w}{\partial u} \\ &= \underbrace{\bar{w}} \underbrace{\frac{\partial w}{\partial u}} \end{aligned}$$

Number which is propagated tras.

Derivative of small local function.

Look up expression of u, v (and/or w) to use.