

Fitting Neural Networks

$$C, \text{ cost } (\underline{f} - \underline{y})^T (\underline{f} - \underline{y}), \text{ or } -\log P(\underline{y} | \underline{f})$$

Prob. model with parameters \underline{f}

train label \underline{y}

$$\underline{f} = g^{(3)}(\underline{a}^{(3)})$$

$$\underline{a}^{(3)} = W^{(3)} \underline{h}^{(2)} + \underline{b}^{(3)}$$

Neural net params:

$$\underline{\theta} = \begin{bmatrix} \text{vec}(W^{(1)}) \\ \text{vec}(\underline{b}^{(1)}) \\ \text{vec}(W^{(2)}) \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$\underline{h}^{(2)} = g^{(2)}(\underline{a}^{(2)})$$

$$\underline{a}^{(2)} = W^{(2)} \underline{h}^{(1)} + \underline{b}^{(2)}$$

$$\underline{h}^{(1)}$$

$$\vdots$$

$$\vdots$$

$$\vdots$$

Training input \underline{x}

SGD fitting:

Initialize $\underline{\theta} \neq \underline{0}$

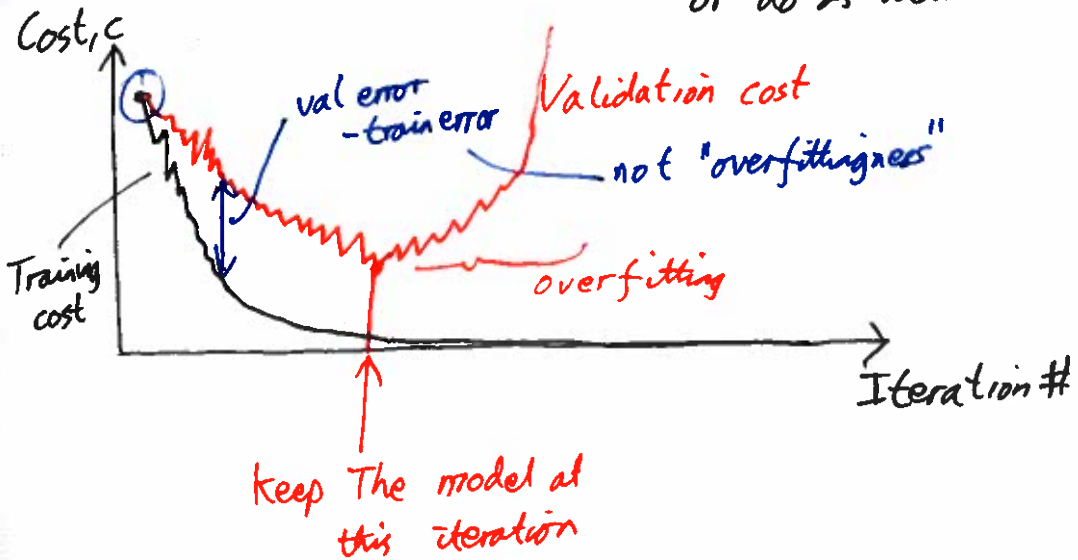
Loop over examples:

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

Today: Termination (early stopping)
Computing $\nabla_{\underline{\theta}} c$ (reverse-mode diff.)

Early Stopping

Form regularization, alternative to L2
or do as well



Every sweep through training set "an epoch":

If val cost is the smallest we've seen

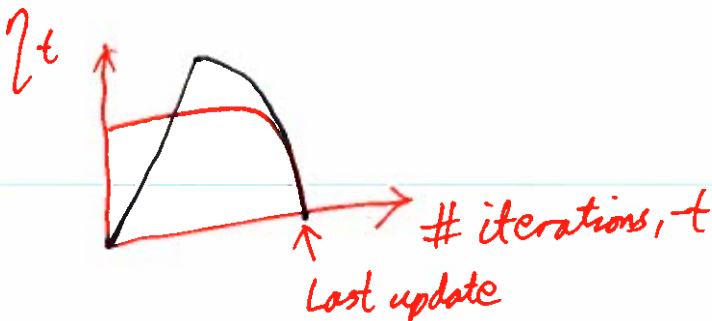
or some # updates

Store weights Θ & val cost

If val. cost hasn't improved in 20 evaluations of val set:

Stop return best Θ

reduce learning rate η



Reverse-mode differentiation "Back-propagation"

Works on a compute graph (DAG)
Directed acyclic

Strategy

For every intermediate Z

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

$$\text{For a matrix } \bar{Z}_{ij} = \frac{\partial c}{\partial Z_{ij}}$$

Start at end of computation

$$\text{Example: } c = (f - y)^2$$

$$\bar{f} = \frac{\partial c}{\partial f} = 2(f - y)$$

$$\text{If have } \underline{f} \text{ and } \underline{y}$$

$k \times 1 \quad k \times 1$

$$c = \sum_k (f_k - y_k)^2$$

$$\bar{f}_j = \frac{\partial c}{\partial f_j} = 2(f_j - y_j)$$

$$\underline{\bar{f}} = 2(\underline{f} - \underline{y})$$

We combine local propagation rules



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

$$\bar{u} = \frac{\partial c}{\partial u} = \underbrace{\frac{\partial c}{\partial w}}_{\bar{w}} \underbrace{\frac{\partial w}{\partial u}}_w$$

Derivative of local function

Depends on u and v and/or w

In example: $\frac{1}{v}$

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

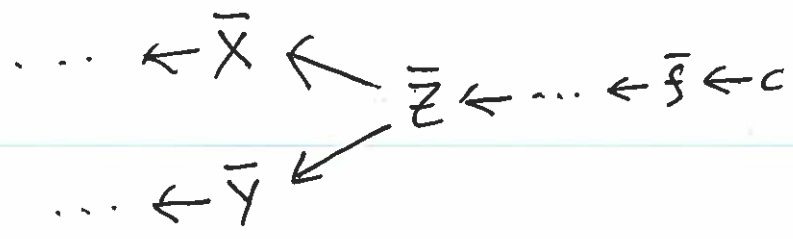
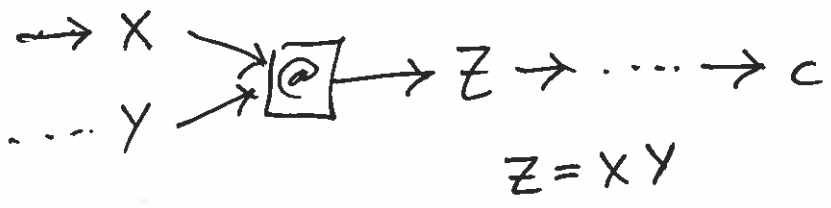
In example $-\frac{u}{v^2}$

Example: Matrix Multiplication

$N \times N$ matrix N^2 numbers in it

$O(\cdot)$ cost of matrix-matrix

N^3 ~~$N^2 \log N$~~ ? ~~$N \log N$~~ ? $N^{2.373}$



$$\bar{X}_{ij} = \frac{\partial c}{\partial X_{ij}} = \sum_{m,n} \underbrace{\frac{\partial c}{\partial Z_{mn}}}_{\bar{Z}_{mn}} \underbrace{\frac{\partial Z_{mn}}{\partial X_{ij}}}_{\delta_{im} Y_{jn}} = \sum_n \bar{Z}_{in} \underbrace{Y_{jn}}_{(Y^T)_{nj}}$$

$$\bar{X} = \bar{Z} Y^T$$

$$\bar{Y} = X^T \bar{Z}$$

$$Z_{mn} = \sum_p X_{mp} Y_{pn}$$

$$= X_{m1} Y_{1n} + X_{m2} Y_{2n} + \dots$$

$$\dots \underbrace{X_{mj} Y_{jn}} + \dots$$

$$\begin{array}{l} Z = X Y \\ M \times N \quad M \times P \quad P \times N \end{array}$$

$$\begin{array}{l} \bar{Y} = X^T \bar{Z} \\ P \times N \quad P \times M \quad M \times N \end{array}$$