

Me: Michael Gelbart

(^{PhD Student}visiting, Mark this year)

mgelbart@seas.harvard.edu

Terminology / Notation

σ : nonlinear element-wise function, "activation function"

\underline{W} : "weights" - free parameters

\underline{b} : "biases" - " "

x : "units", "neurons", "hidden units"

train / learn / optimize \rightarrow set W, b

Supervised learning: learn $y=f(x)$
from labeled data

$$f(x) = 2x$$

function

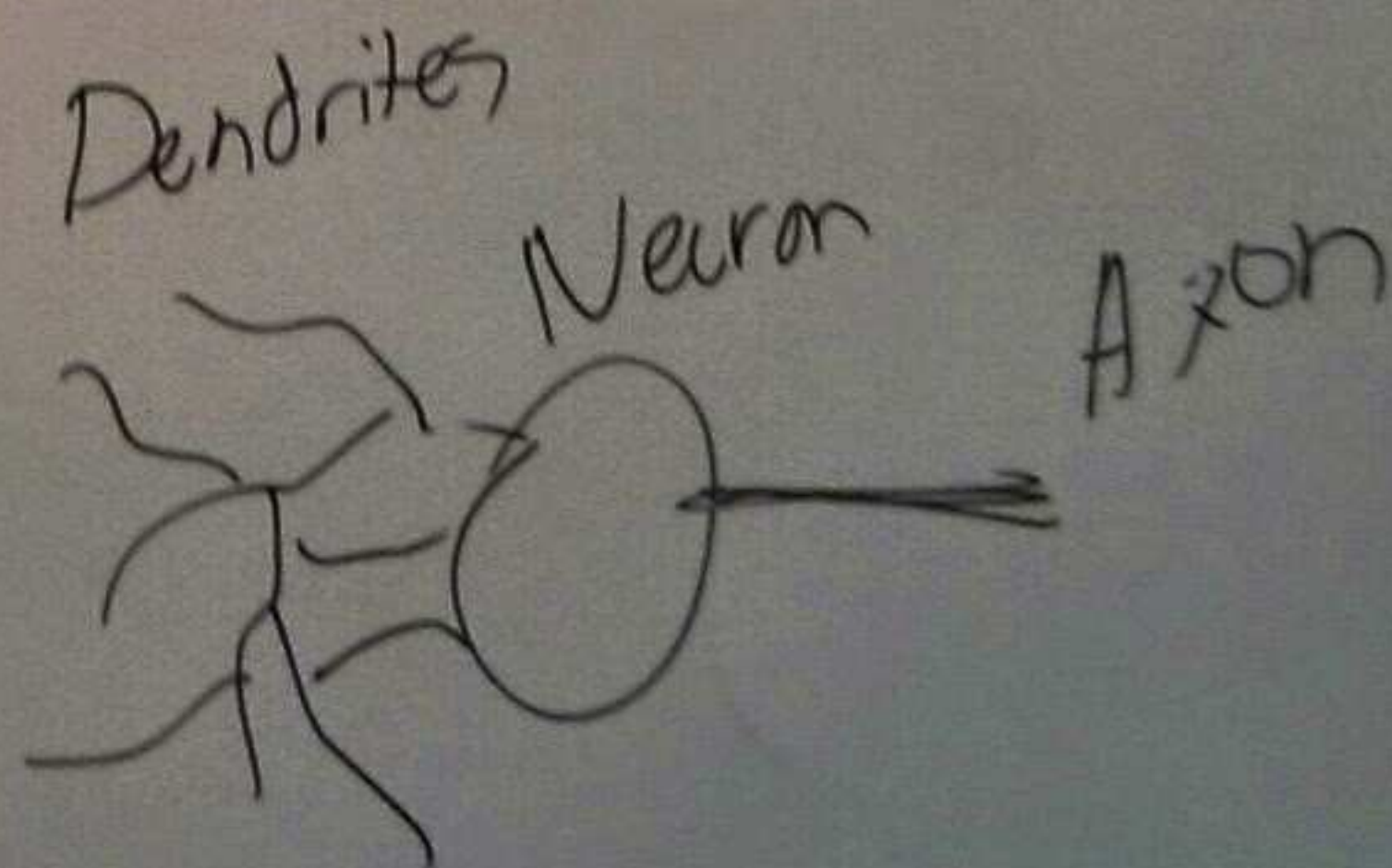
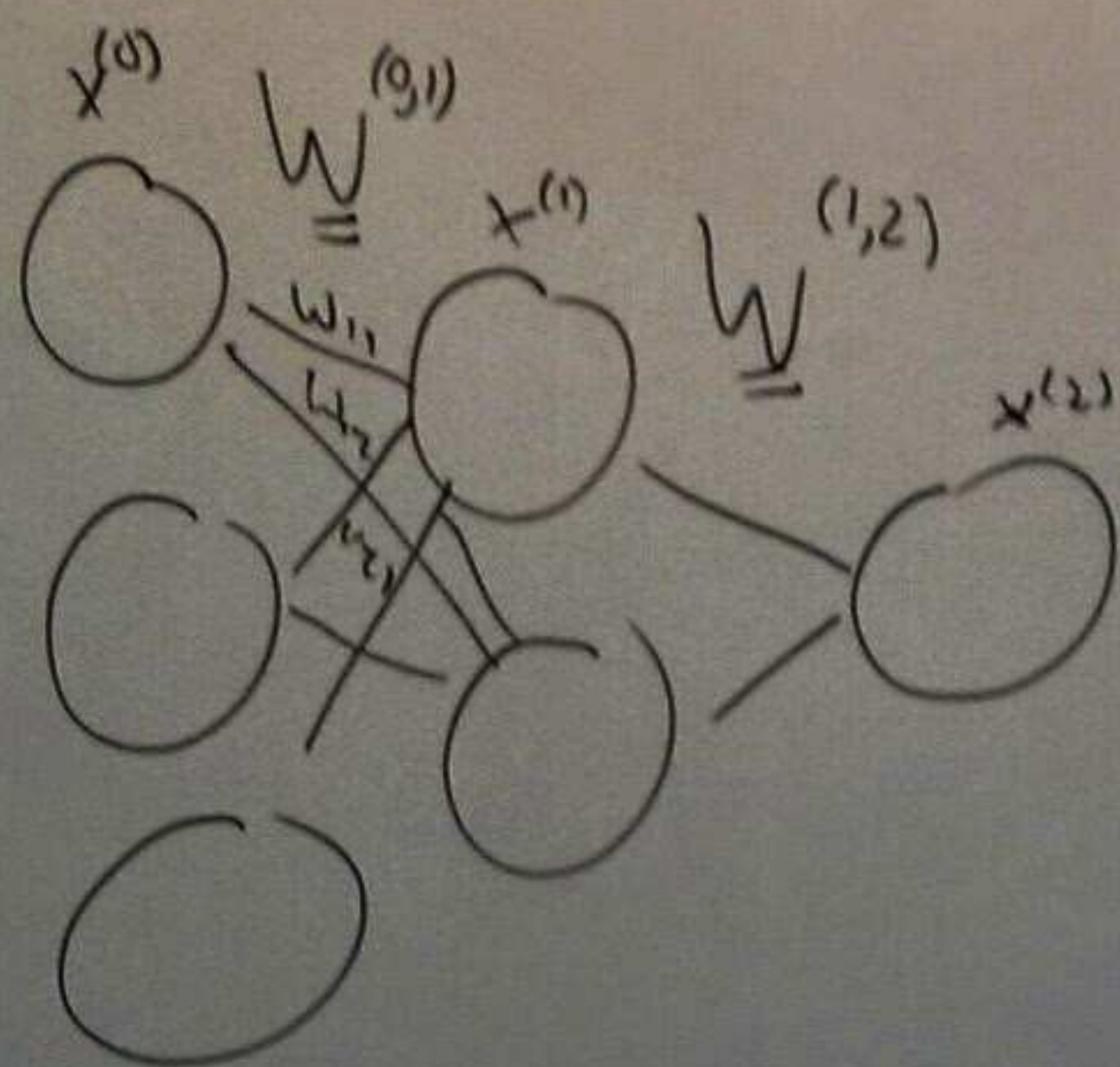
$$f(x) = ax^b$$

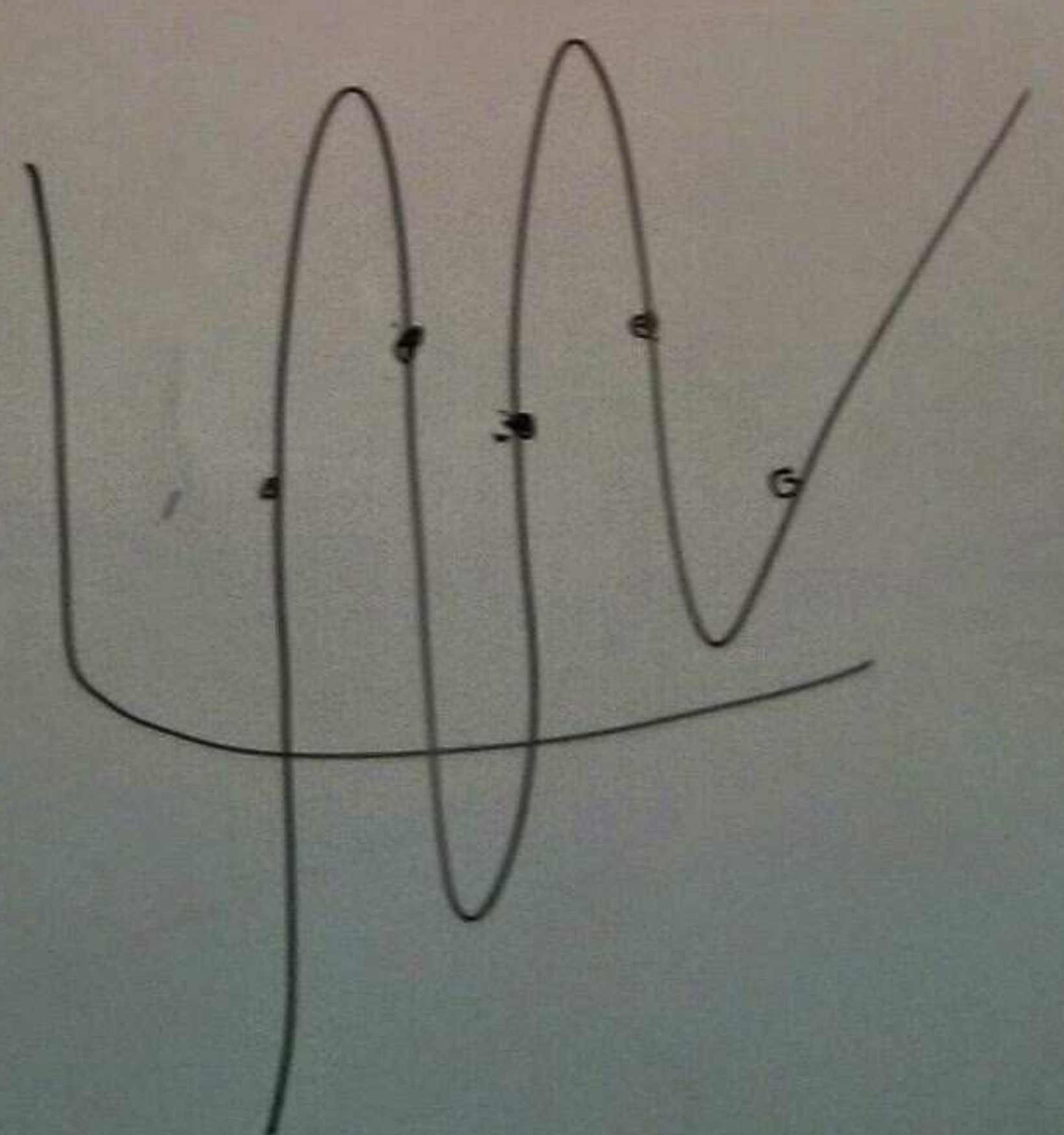
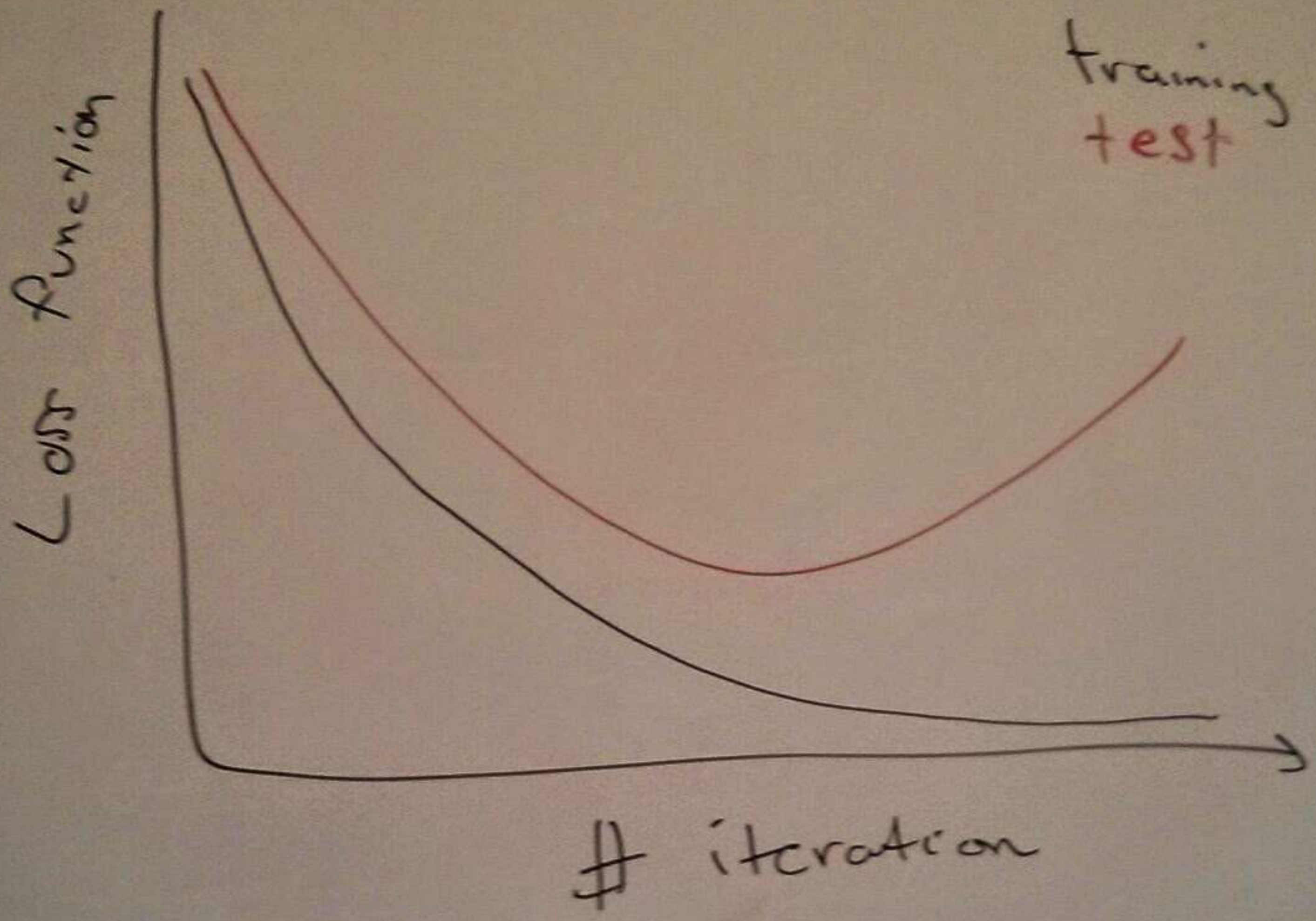
class of function

polynomials

Neural networks

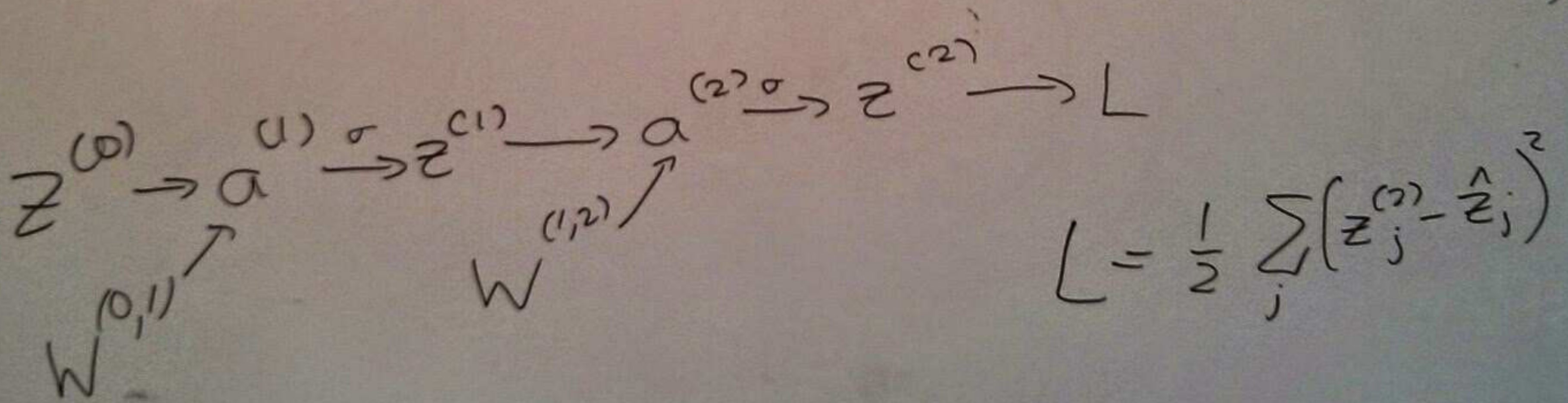
$$X^{(j+1)} = \sigma \left(\sum_i W_{ij}^{(j+1)} X_i^{(j)} + b^{(j)} \right)$$





Backpropagation

is an ^{efficient} algorithm to compute the gradient of the loss w.r.t. weights (making extensive use of the chain rule)



$$z_j^{(n)} = \sigma(a_j^{(n)})$$

$$a_j^{(n)} = \sum_i W_{ji}^{(n-1,n)} z_i^{(n-1)}$$

$$z_i^{(1)}$$



$$\frac{\partial L}{\partial W_{ji}^{(1,2)}} = \frac{\partial L}{\partial a_j^{(2)}} \frac{\partial a_j^{(2)}}{\partial W_{ji}^{(1,2)}}$$

$$\frac{\partial L}{\partial a_j^{(2)}} = \frac{\partial L}{\partial z_j^{(2)}} \frac{\partial z_j^{(2)}}{\partial a_j^{(2)}}$$

$$z_j^{(2)} - \hat{z}_j$$

$$\sigma'(a_j^{(2)})$$

$$\frac{\partial L}{\partial W_{ji}^{(1,2)}} = (z_j^{(2)} - \hat{z}_j) \sigma'(a_j^{(2)}) z_i^{(1)}$$

$$\frac{\partial L}{\partial W_{ji}^{(0,1)}} = \frac{\partial L}{\partial a_j^{(1)}} \frac{\partial a_j^{(1)}}{\partial W_{ji}^{(0,1)}} = z_i^{(0)} \frac{\partial L}{\partial z_j^{(1)}} \frac{\partial z_j^{(1)}}{\partial a_j^{(1)}}$$

$$= z_i^{(0)} \sigma'(a_j^{(1)}) \frac{\partial L}{\partial z_j^{(1)}}$$

$$= z_i^{(0)} \sigma'(a_j^{(1)}) \sum_k \frac{\partial a_k}{\partial z_j} \frac{\partial L}{\partial a_k}$$

$$= z_i^{(0)} \sigma'(a_j^{(1)}) \sum_k w_{kj}^{(1,2)} (z_k^{(2)} - \hat{z}_k) \sigma'(a_k^{(2)})$$