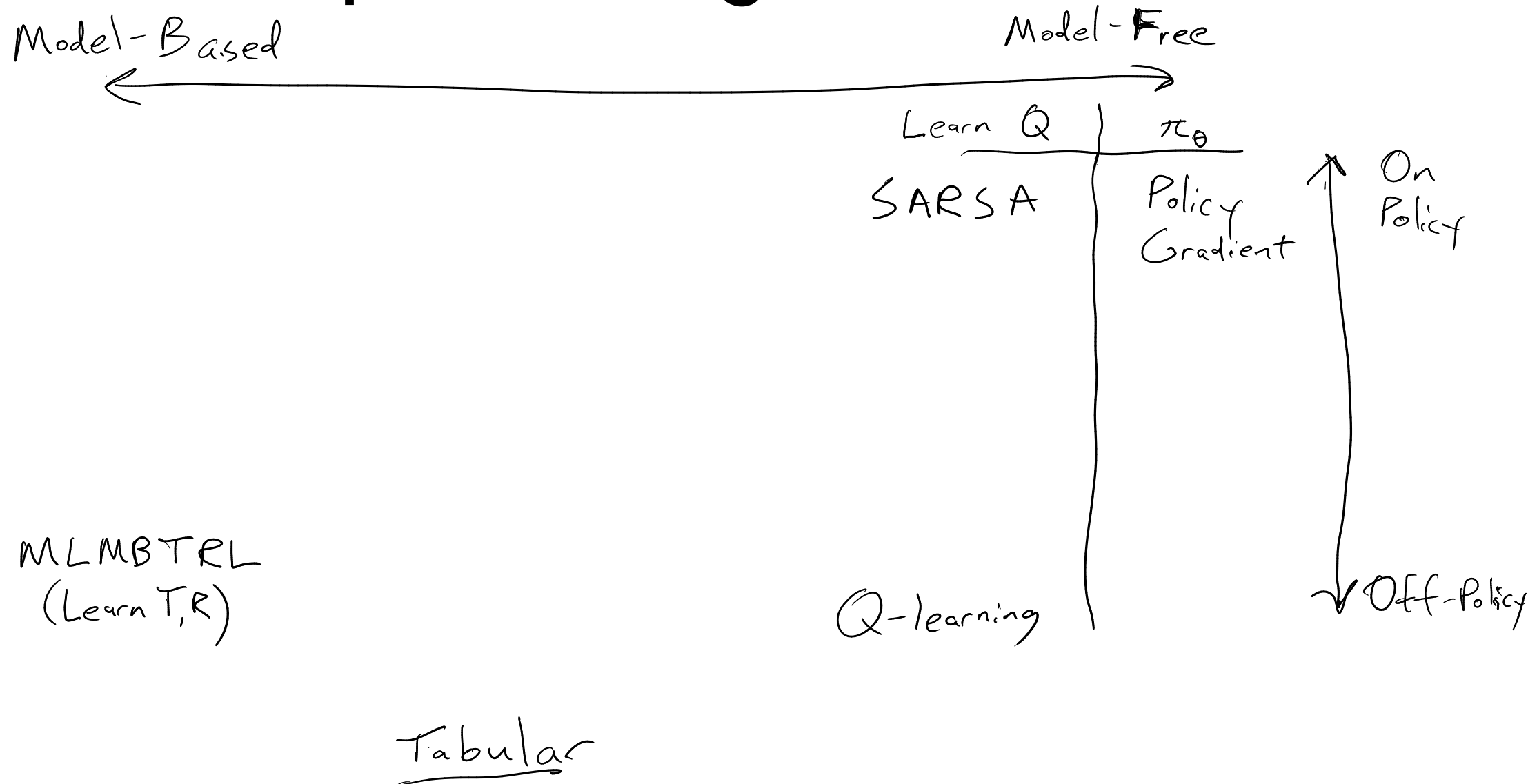


Neural Network Function Approximation

Map of RL Algorithms



This Time

Challenges in Reinforcement Learning:


- Exploration vs Exploitation
- Credit Assignment
- Generalization

Function Approximation

Function Approximation

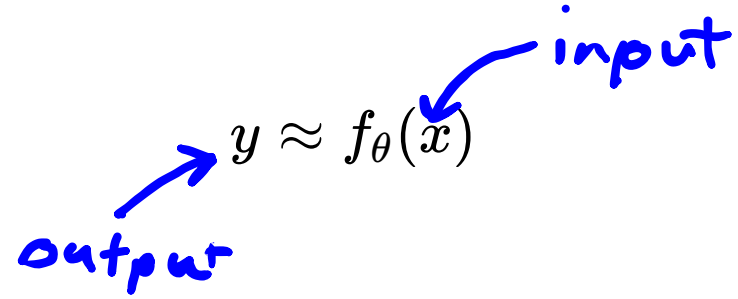
$$y \approx f_{\theta}(x)$$

Function Approximation

$$y \approx f_{\theta}(x)$$


input

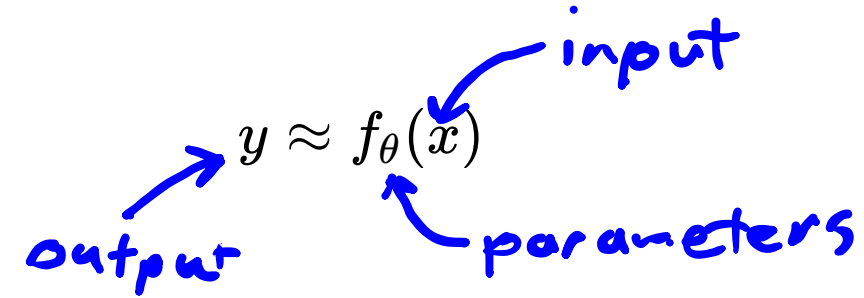
Function Approximation



A diagram illustrating the function approximation equation $y \approx f_{\theta}(x)$. The equation is written in black. A blue arrow points from the handwritten word "output" to the variable y . Another blue arrow points from the handwritten word "input" to the variable x inside the function notation.

$$y \approx f_{\theta}(x)$$

Function Approximation



A diagram illustrating the function approximation equation $y \approx f_{\theta}(x)$. The equation is centered, with three handwritten blue arrows pointing to its components: one from the word "output" to y , one from the word "input" to x , and one from the word "parameters" to θ .

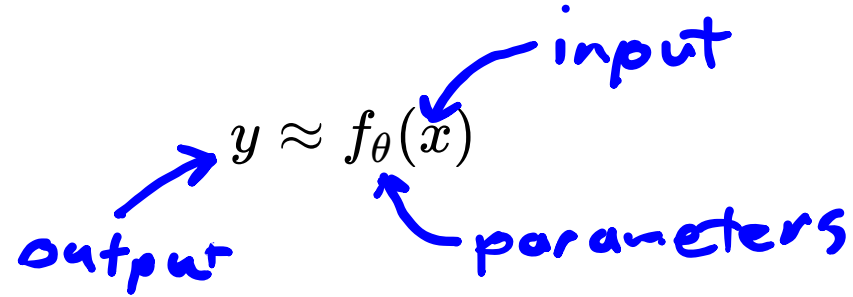
$$y \approx f_{\theta}(x)$$

output

input

parameters

Function Approximation



A diagram illustrating the function approximation equation $y \approx f_{\theta}(x)$. The equation is written in black. Three blue arrows point to different parts of the equation: one from the word "output" to y , one from the word "input" to x , and one from the word "parameters" to θ .

Previously, Linear:

$$f_{\theta}(x) = \theta^{\top} \beta(x)$$

Function Approximation

$$y \approx f_{\theta}(x)$$

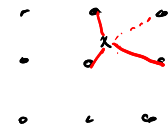
output input parameters

Previously, Linear:

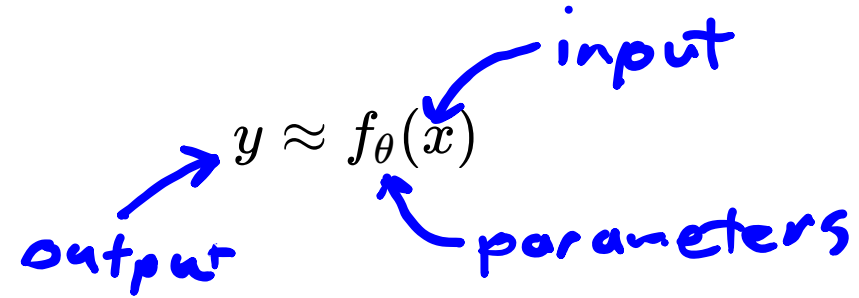
$$f_{\theta}(x) = \theta^{\top} \beta(x)$$

weights features

interpolation
weight

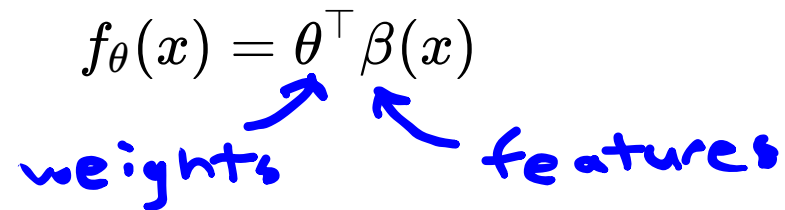


Function Approximation



A diagram showing the equation $y \approx f_{\theta}(x)$. A blue arrow points from the handwritten word "output" to the variable y . Another blue arrow points from the handwritten word "input" to the variable x . A third blue arrow points from the handwritten word "parameters" to the symbol θ .

Previously, Linear:



A diagram showing the equation $f_{\theta}(x) = \theta^{\top} \beta(x)$. A blue arrow points from the handwritten word "weights" to the symbol θ . Another blue arrow points from the handwritten word "features" to the symbol β .

e.g. $\beta_i(x) = \sin(i \pi x)$

AI = Neural Nets

**Neural Nets are
just another
function
approximator**

Neural Network

Neural Network

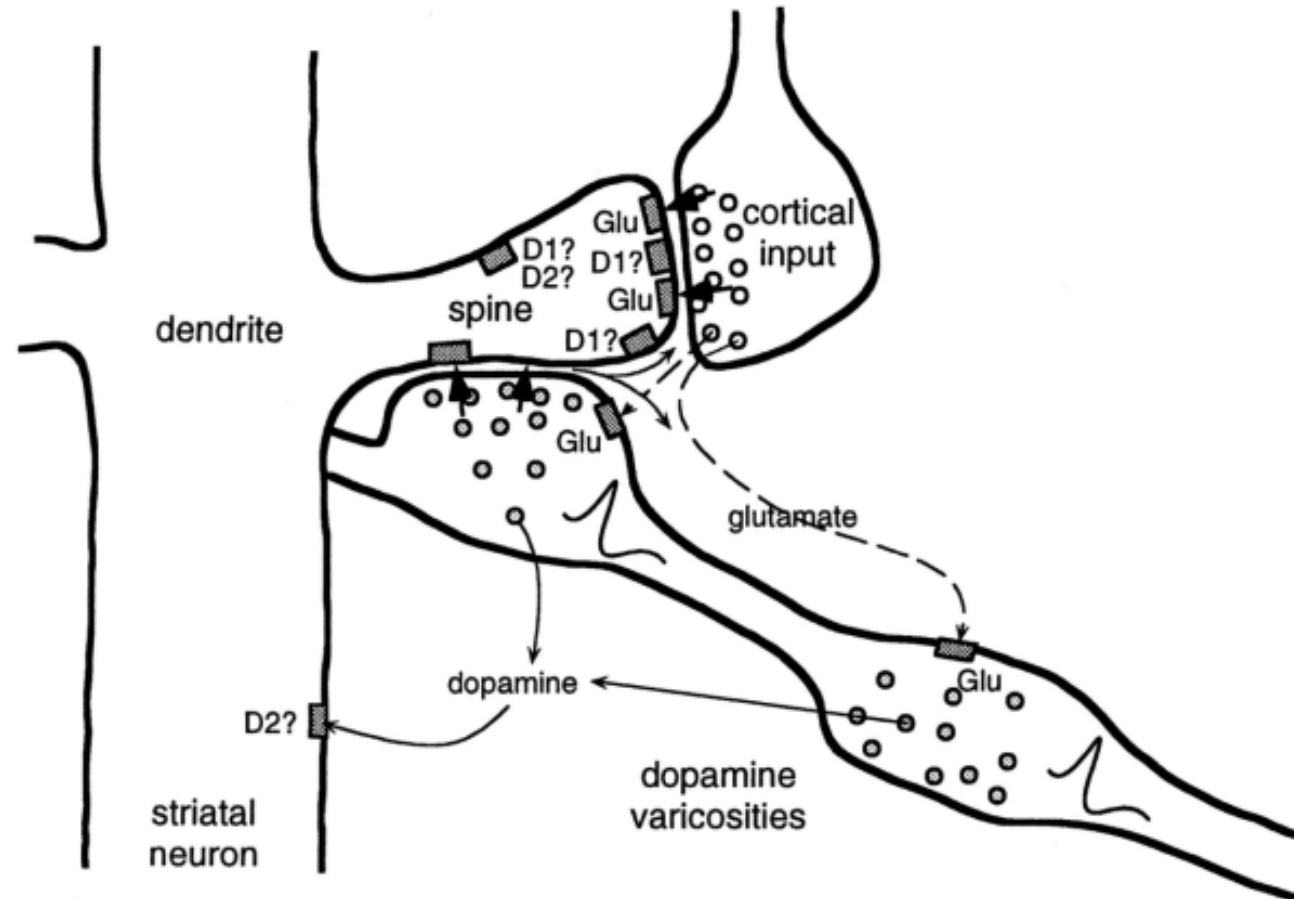
$$h(x) = \sigma(Wx + b)$$

nonlinearity

$$w^{(2)} \sigma(w^{(1)}x + b^{(1)}) + b^{(2)}$$

Neural Network

$$h(x) = \sigma(Wx + b)$$



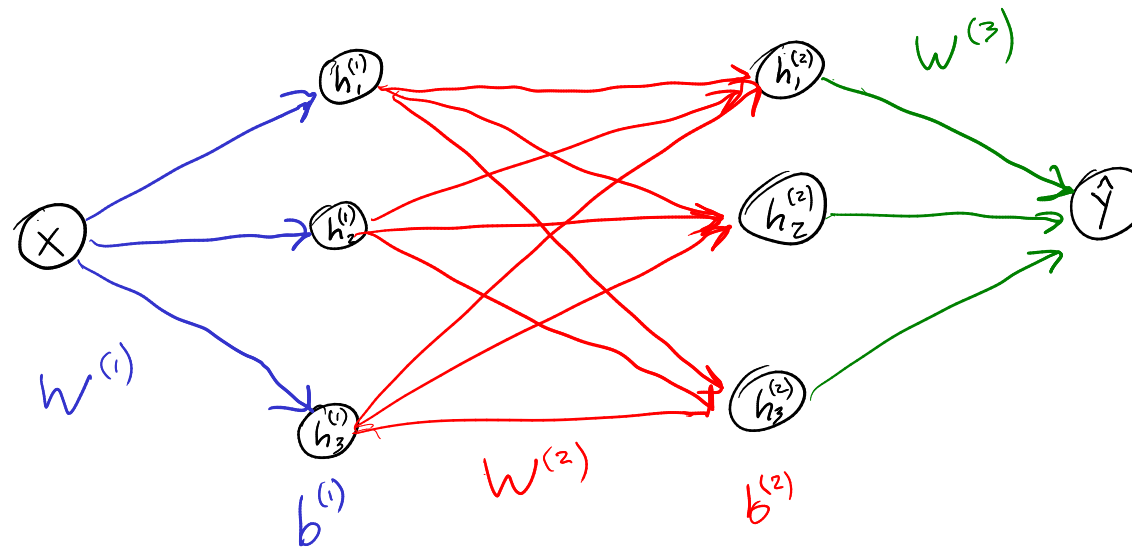
Neural Network

Neural Network

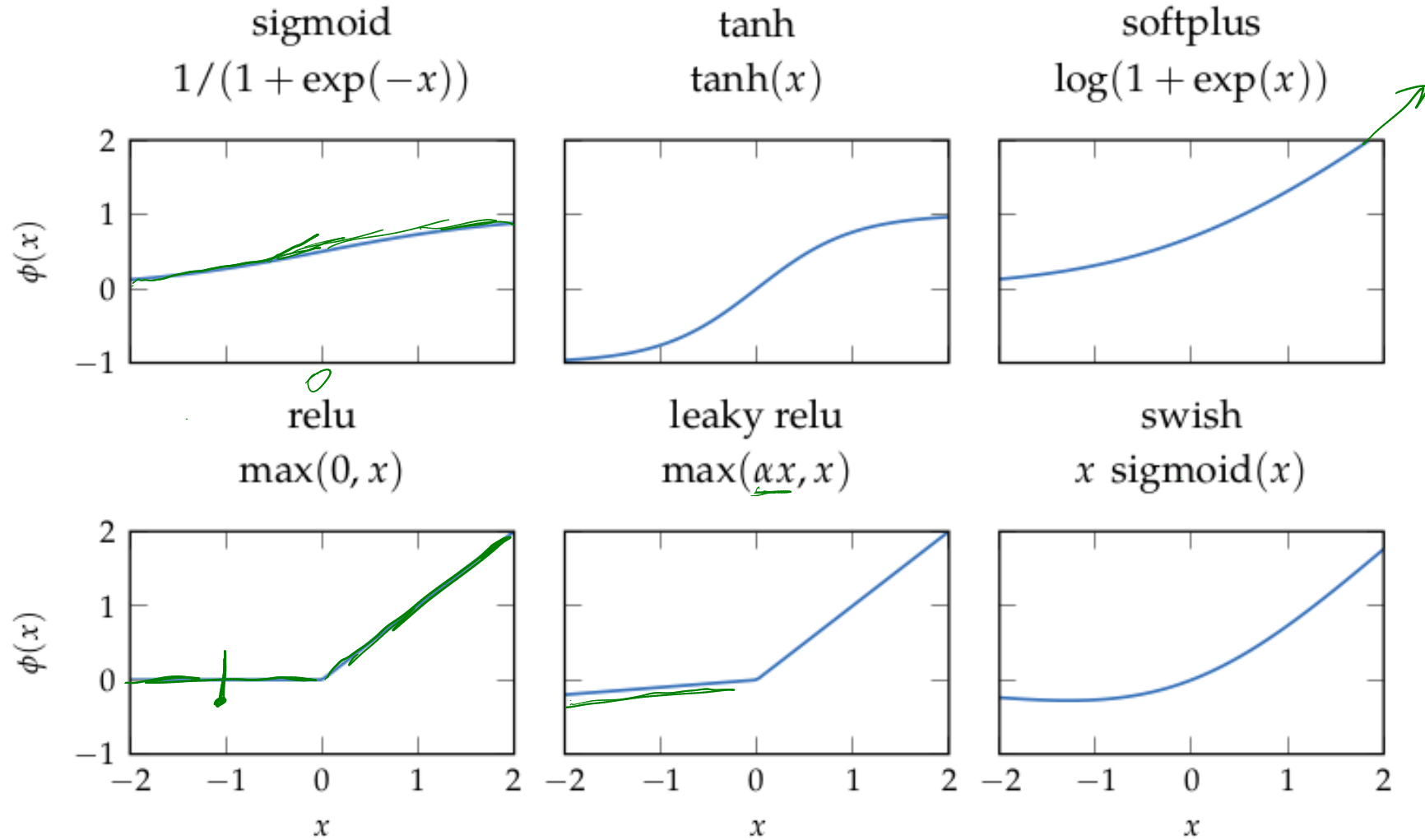
$$\theta = (w^{(3)}, w^{(2)}, w^{(1)}, b^{(2)}, b^{(1)})$$

$$h(x) = \sigma(Wx + b)$$

$$f_{\theta}(x) = h^{(2)}(h^{(1)}(x)) = \underline{w^{(3)}} \sigma^{(2)}(\underline{w^{(2)}} \sigma^{(1)}(\underline{w^{(1)}} x + \underline{b^{(1)}}) + \underline{b^{(2)}})$$

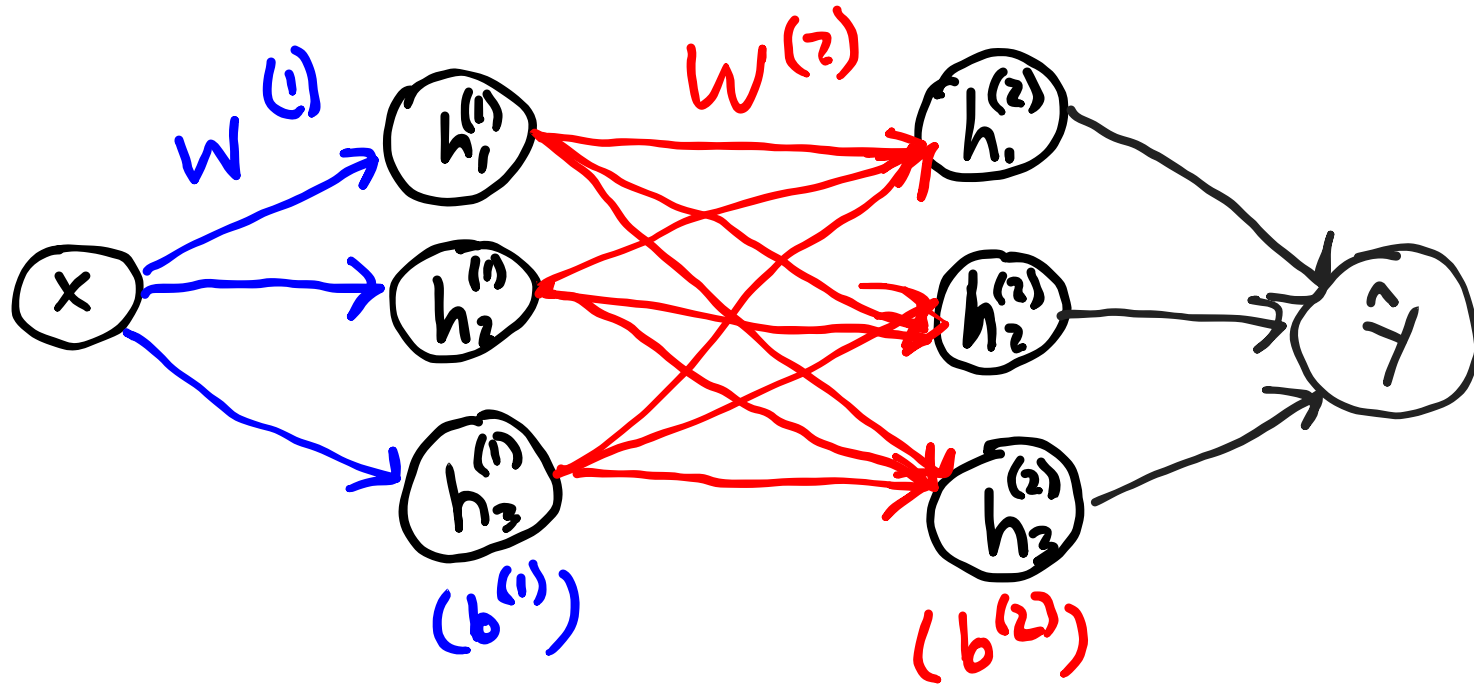


Nonlinearities

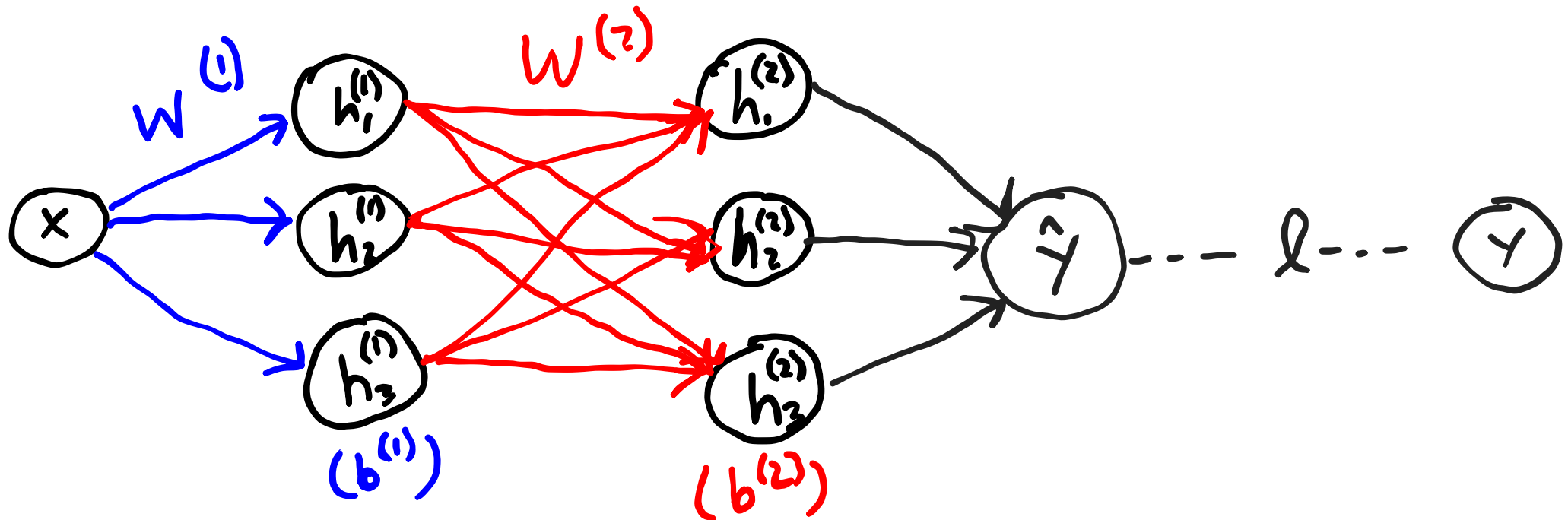


Training

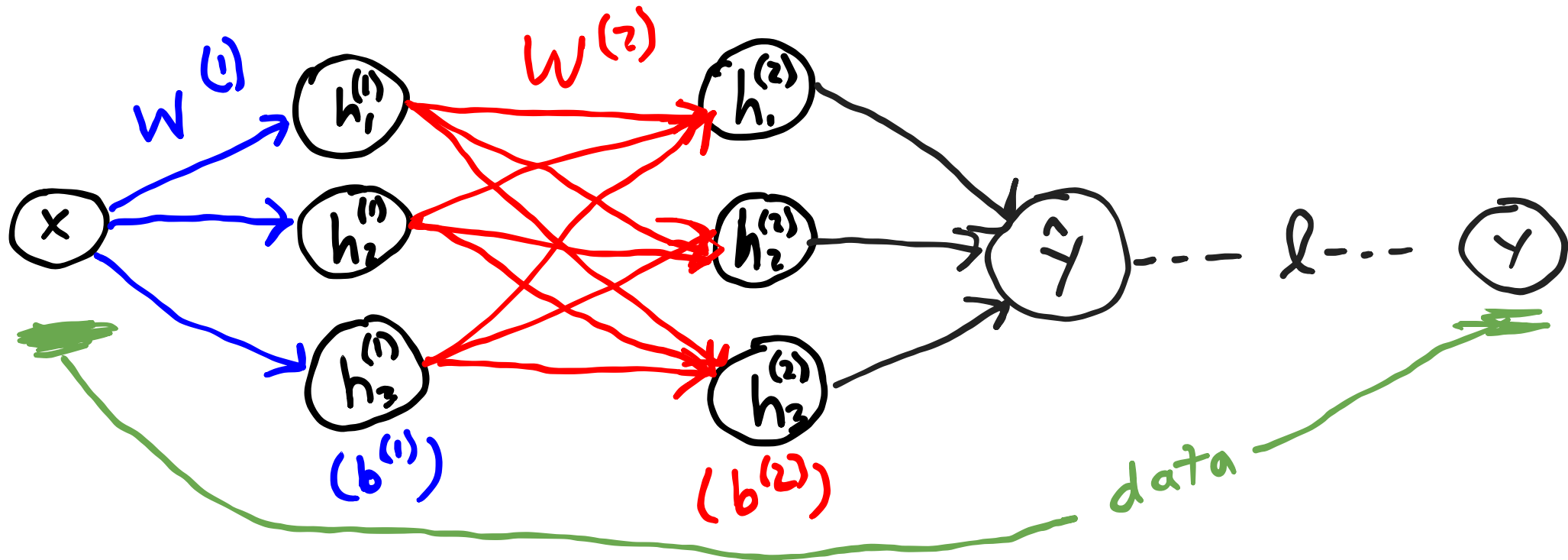
Training



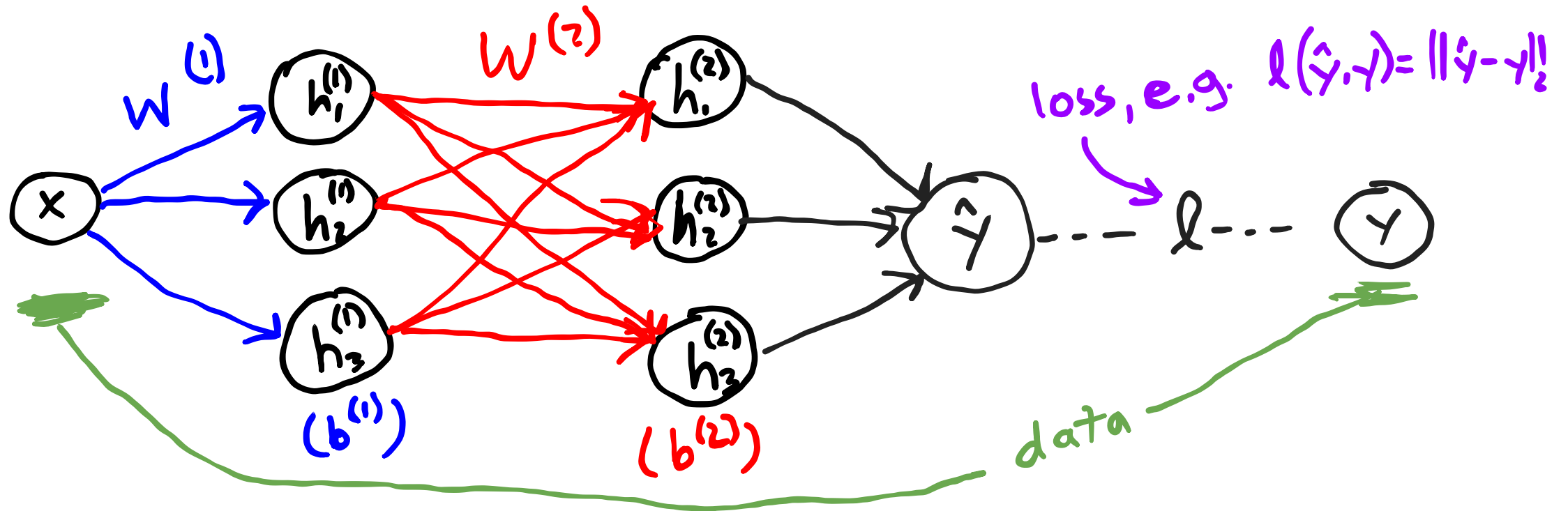
Training



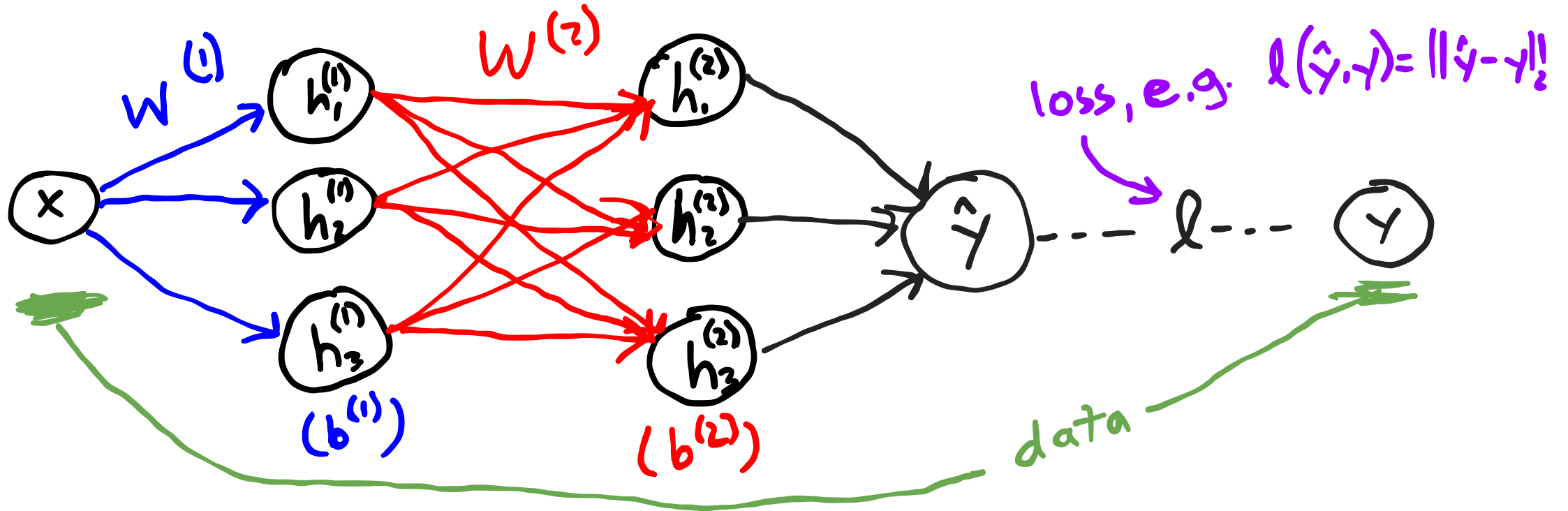
Training



Training

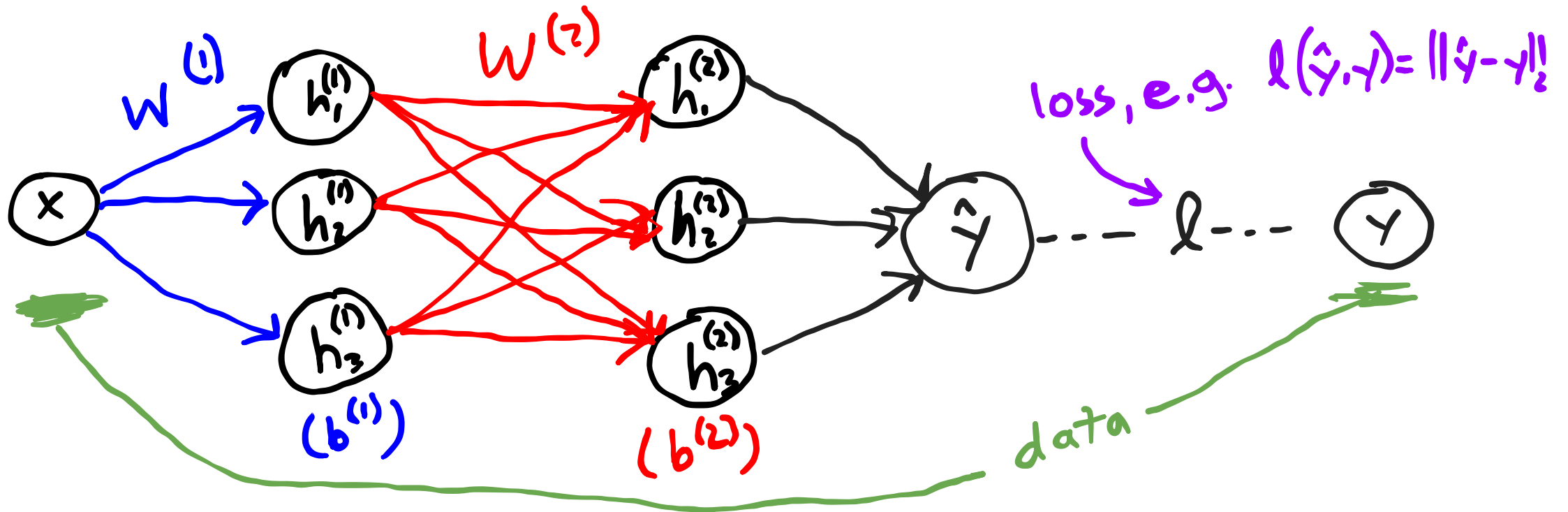


Training



$$\theta^* = \arg \min_{\theta} \sum_{(x,y) \in \mathcal{D}} l(f_{\theta}(x), y)$$

Training



$$\theta^* = \arg \min_{\theta} \sum_{(x,y) \in \mathcal{D}} l(f_{\theta}(x), y)$$

Stochastic Gradient Descent: $\theta \leftarrow \theta - \alpha \nabla_{\theta} l(f_{\theta}(x), y)$

Chain Rule

$$f \circ g \circ h = f(g(h(x)))$$

$$\nabla_{\theta} l(f_{\theta}(x), y)$$

$$\theta = (w^{(1)}, w^{(2)}, b^{(1)}, b^{(2)})$$

$$l(\hat{y}, y) = \frac{1}{n} \sum_{i=1}^n (\hat{y} - y)^2$$

$$\frac{\partial l}{\partial \hat{y}} = -\frac{1}{n} 2(\hat{y} - y)$$

$$\left. \frac{\partial f(g(h(x)))}{\partial x} \right|_{x_0} = \left. \frac{\partial f(g(h))}{\partial h} \right|_{h_0} \left. \frac{\partial h(x)}{\partial x} \right|_{x_0} = \left. \frac{\partial f(g)}{\partial g} \right|_{g_0} \left. \frac{\partial g(h)}{\partial h} \right|_{h_0} \left. \frac{\partial h(x)}{\partial x} \right|_{x_0}$$

$$\hat{y} = \overbrace{w^{(2)}}^{f_{\theta}(x)} \sigma(\underbrace{w^{(1)} x + b^{(1)}}_{\text{blue underline}}) + b^{(2)}$$

$$\left. \frac{\partial l}{\partial w^{(2)}} \right|_{x_0} = \left. \frac{\partial l}{\partial \hat{y}} \right|_{\hat{y}_0} \left(\left. \frac{\partial \hat{y}}{\partial w^{(2)}} \right|_{x_0} \right) = \underbrace{\left. \frac{\partial l}{\partial \hat{y}} \right|_{\hat{y}_0}}_{\text{red underline}} \underbrace{\left(\sigma(w^{(1)} x_0 + b^{(1)}) \right)}_{\text{blue underline}}$$

SGD

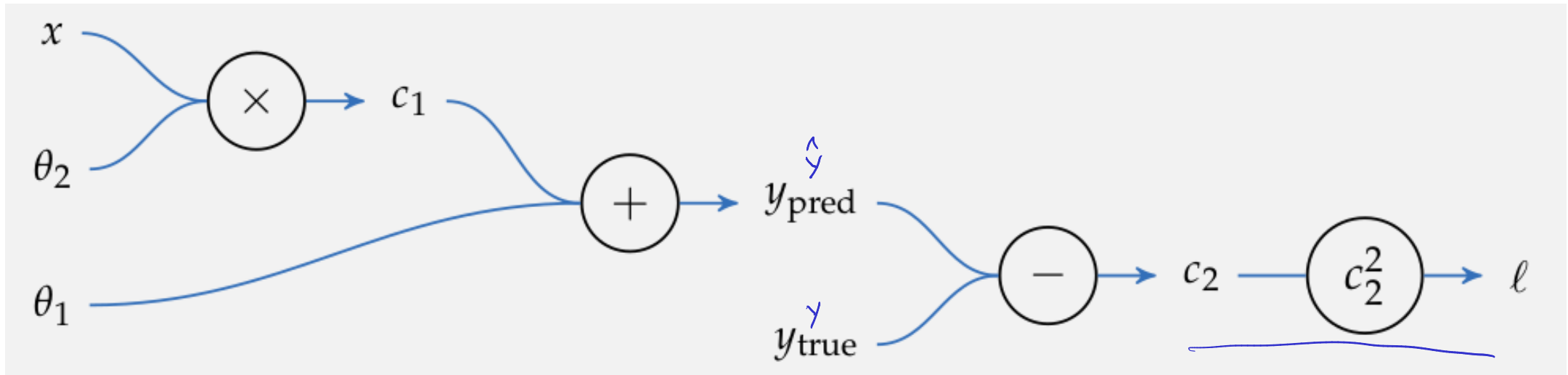
$$w^{(2)} \leftarrow w^{(2)} - \alpha \frac{\partial l}{\partial w^{(2)}}$$

Backprop

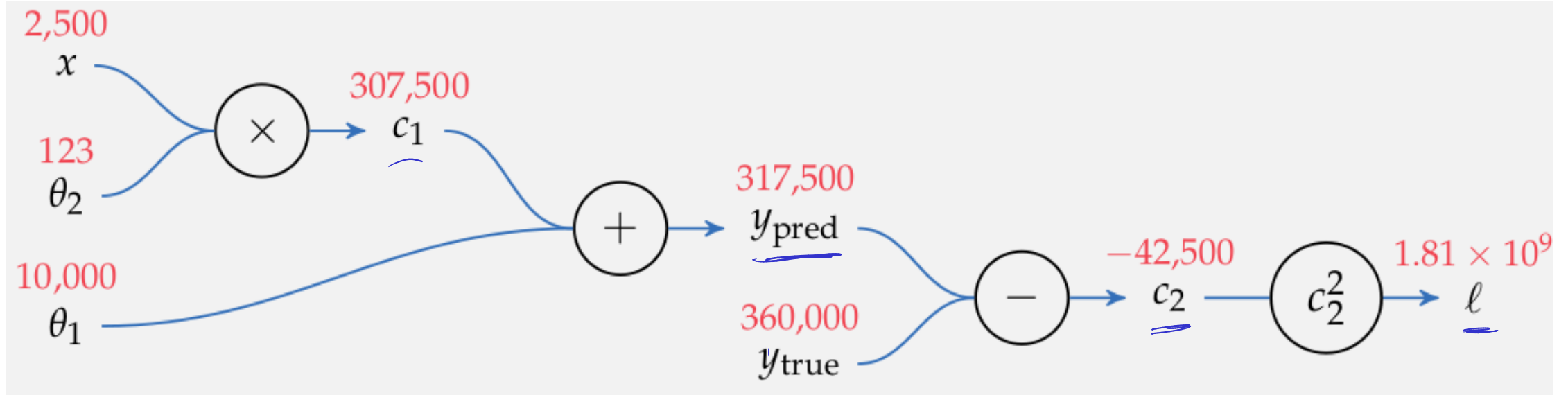
Backprop

$$\hat{y} = \theta_2 x + \theta_1$$

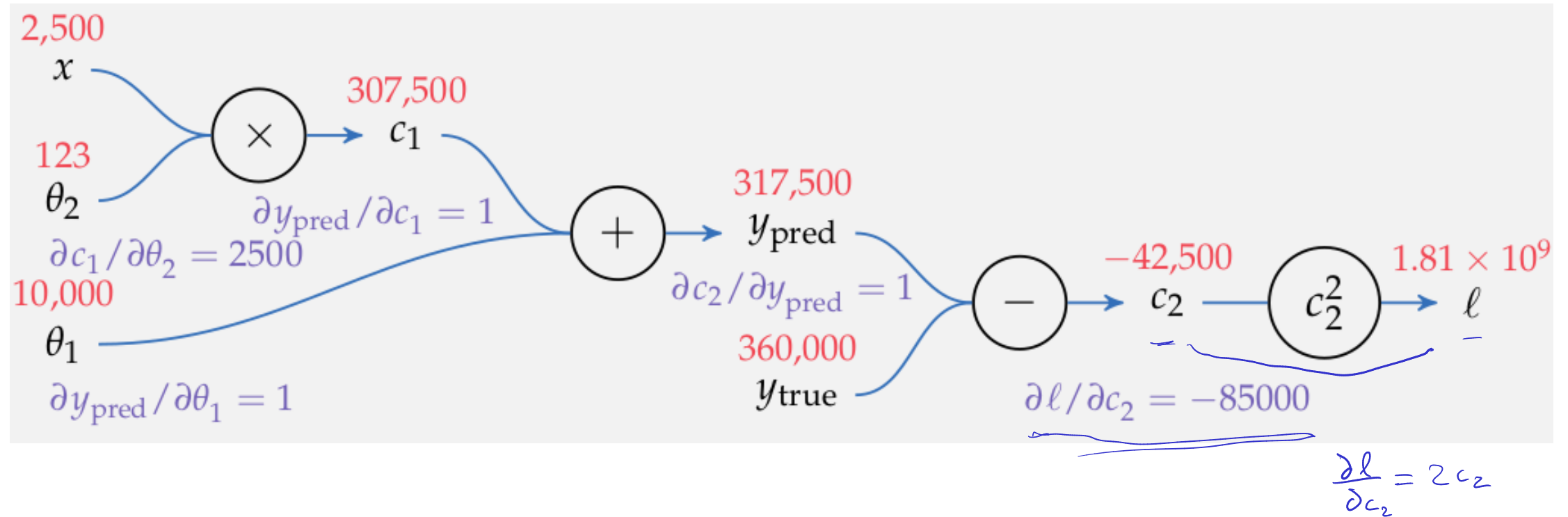
$$l(\hat{y}, y) = (\hat{y} - y)^2$$



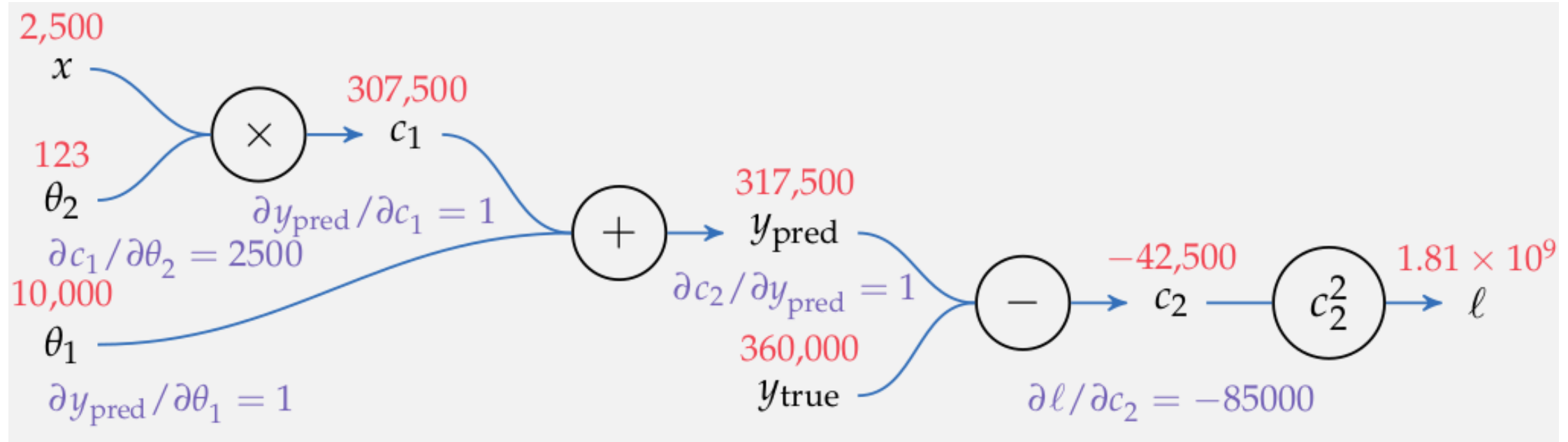
Backprop



Backprop



Backprop



$$\frac{\partial \ell}{\partial \theta_1} = \frac{\partial \ell}{\partial c_2} \frac{\partial c_2}{\partial y_{\text{pred}}} \frac{\partial y_{\text{pred}}}{\partial \theta_1} = -85,000 \cdot 1 \cdot 1 = -85,000$$

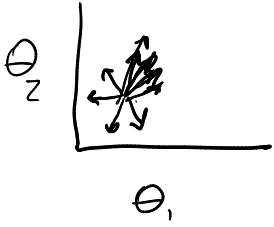
$$\frac{\partial \ell}{\partial \theta_2} = \frac{\partial \ell}{\partial c_2} \frac{\partial c_2}{\partial y_{\text{pred}}} \frac{\partial y_{\text{pred}}}{\partial c_1} \frac{\partial c_1}{\partial \theta_2} = -85,000 \cdot 1 \cdot 1 \cdot 2,500 = -2.125 \times 10^8$$

a “fast and furious” approach to training neural networks does not work and only leads to suffering. Now, suffering is a perfectly natural part of getting a neural network to work well, but it can be mitigated by being thorough, defensive, paranoid, and obsessed with visualizations of basically every possible thing. The qualities that in my experience correlate most strongly to success in deep learning are patience and attention to detail.

Keep calm and lower
your learning rate

- Andrej Karpathy

Adaptive Step Size: RMSProp



$$\Theta \leftarrow \Theta + \alpha \underbrace{\nabla_{\Theta} f_{\Theta}(x)}_{g^{(k)}}$$



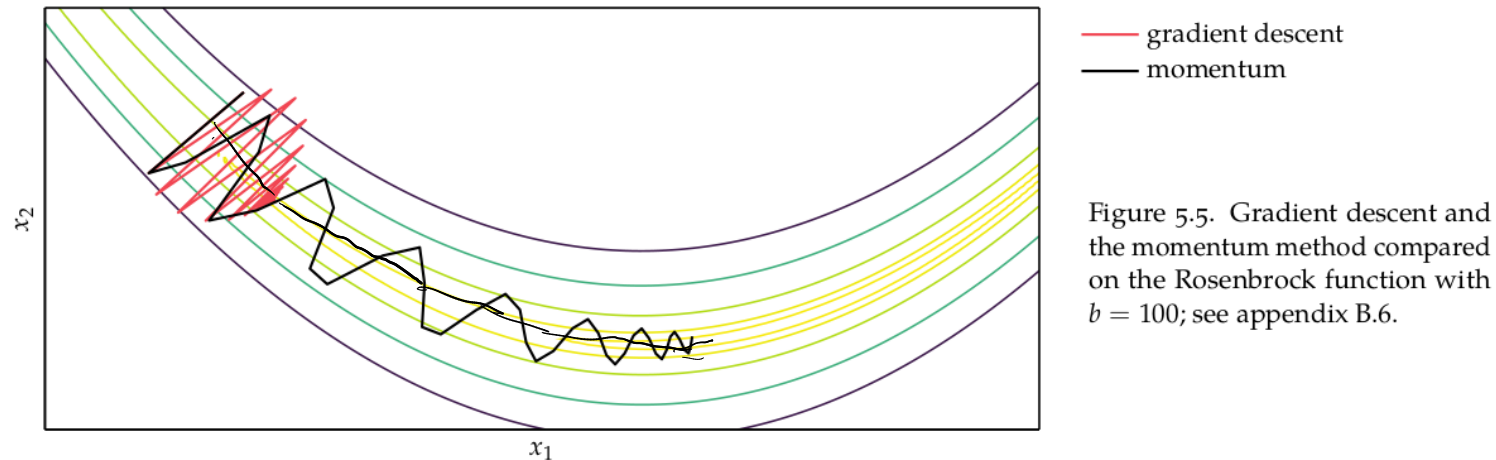
$$\hat{s}^{(k+1)} = \gamma \hat{s}^{(k)} + (1-\gamma) (g^{(k)} \odot g^{(k)})$$

↑ element-wise product

$$\theta_i^{(k+1)} = \theta_i^{(k)} - \frac{\alpha}{\epsilon + \sqrt{\hat{s}_i^{(k+1)}}} g_i^{(k)}$$

Adaptive Step Size: ADAM

(Adaptive Moment Estimation)



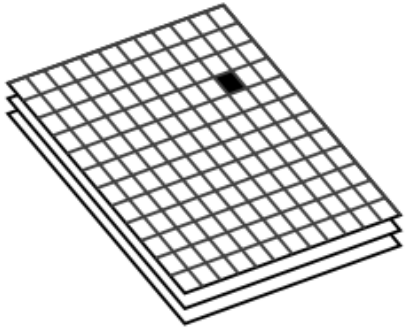
biased decaying momentum $v^{(k+1)} = \gamma_v v^{(k)} + (1 - \gamma_v) g^{(k)}$

biased decaying sq. grad. $s^{(k+1)} = \gamma_s s^{(k)} + (1 - \gamma_s) (g^{(k)} \odot g^{(k)})$

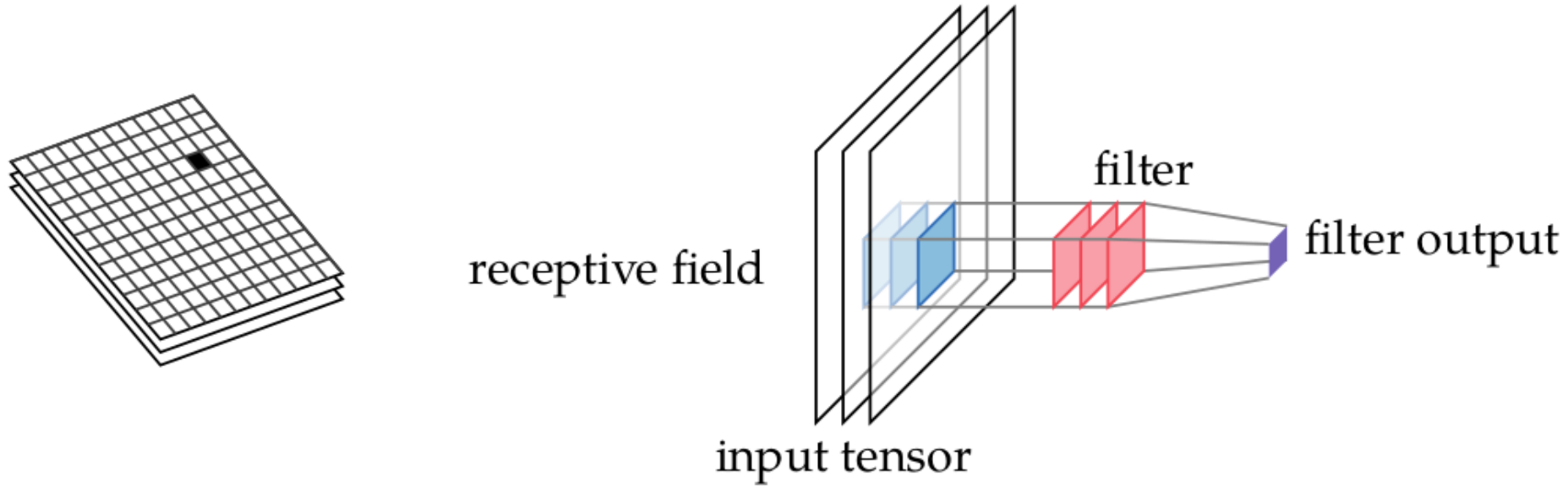
corrected decaying momentum \hat{v}

On Your Radar: ConvNets

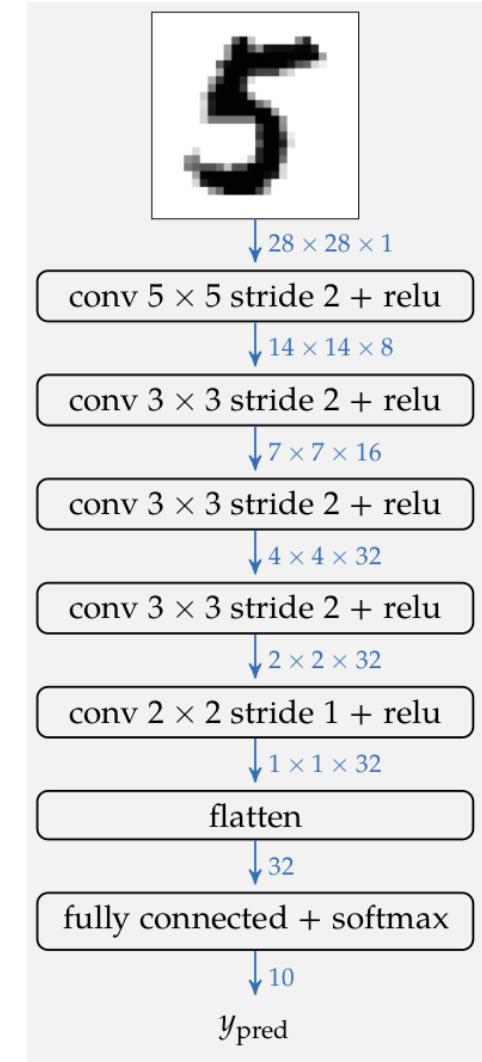
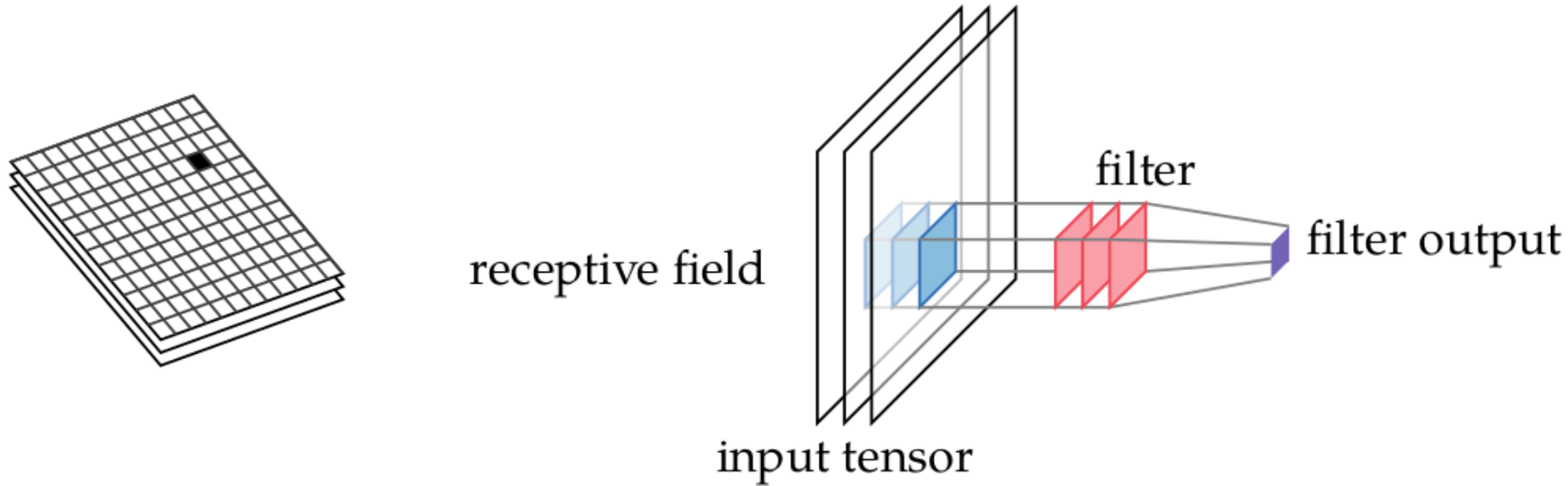
On Your Radar: ConvNets



On Your Radar: ConvNets



On Your Radar: ConvNets



On Your Radar: Regularization

On Your Radar: Regularization

$$\arg \min_{\boldsymbol{\theta}} \sum_{(x,y) \in \mathbf{D}} \ell(f_{\boldsymbol{\theta}}(x), y) - \beta \|\boldsymbol{\theta}\|^2$$

On Your Radar: Regularization

$$\arg \min_{\theta} \sum_{(x,y) \in \mathbf{D}} \ell(f_{\theta}(x), y) - \beta \|\theta\|^2$$

e.g. Batch norm, layer norm, dropout

On Your Radar: Skip Connections (Resnets)

Resources

OpenAI Spinning up