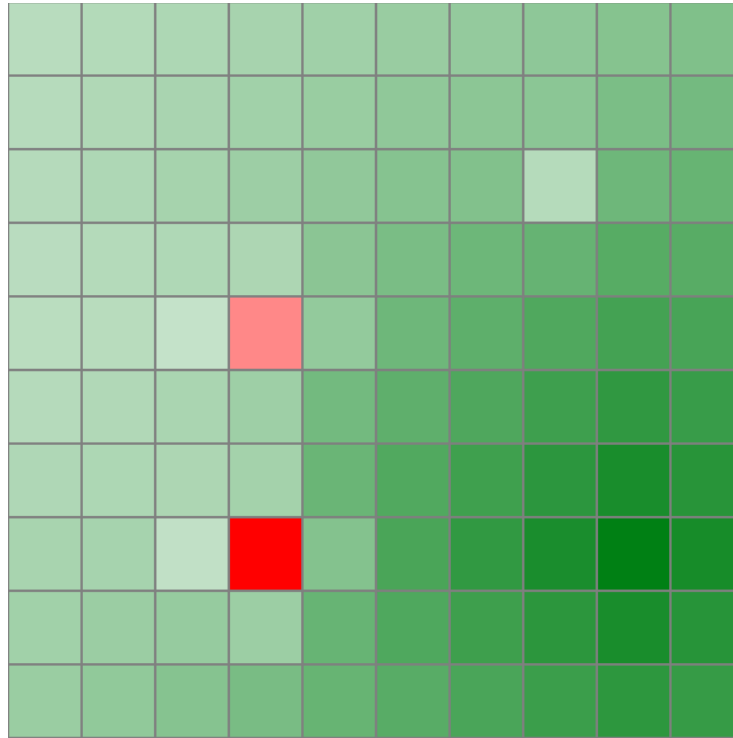


# Neural Network Function Approximation




Do we really need to keep track of  $U(s)$  for every  $U$  separately?

# Function Approximation

# Function Approximation

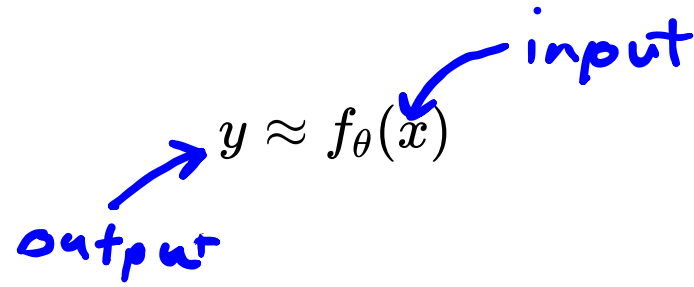
$$y \approx f_{\theta}(x)$$

# Function Approximation

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


input

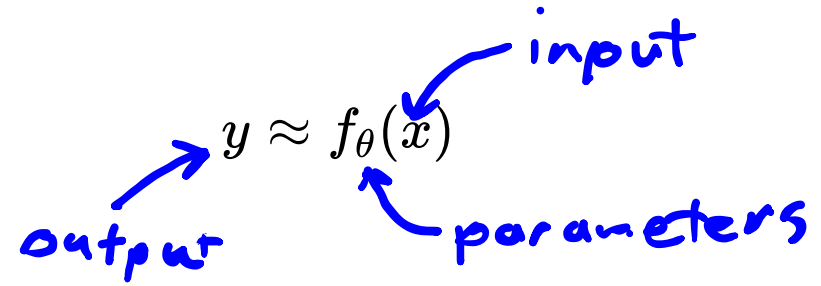
# Function Approximation



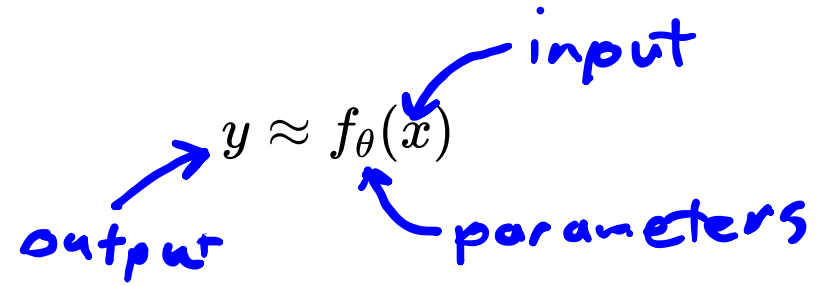
A diagram illustrating the function approximation equation  $y \approx f_{\theta}(x)$ . The equation is centered, with a handwritten blue arrow pointing from the word "output" to the variable  $y$  on the left, and another handwritten blue arrow pointing from the word "input" to the variable  $x$  inside the function  $f_{\theta}(x)$  on the right.

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

# Function Approximation



# Function Approximation

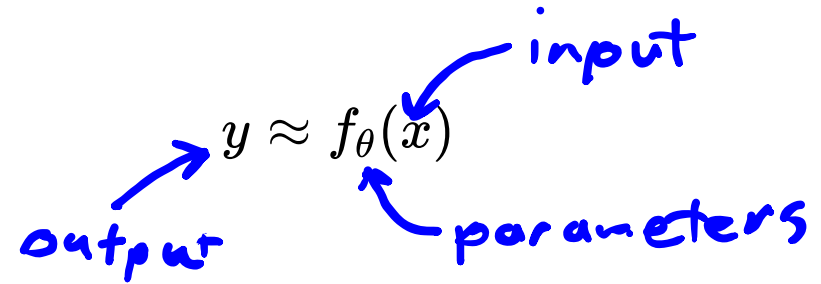


Example: Linear Function Approximation:

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



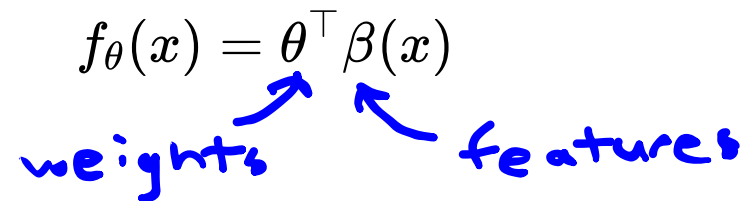
# Function Approximation



A diagram showing the equation  $y \approx f_{\theta}(x)$ . Three blue arrows point to the components: one from the word "output" to  $y$ , one from the word "input" to  $x$ , and one from the word "parameters" to  $\theta$ .

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

Example: Linear Function Approximation:



A diagram showing the equation  $f_{\theta}(x) = \theta^{\top} \beta(x)$ . Two blue arrows point to the components: one from the word "weights" to  $\theta$ , and one from the word "features" to  $\beta(x)$ .

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

# Function Approximation

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

Diagram illustrating the function approximation equation  $y \approx f_{\theta}(x)$  with handwritten blue annotations:

- An arrow points from the word "output" to  $y$ .
- An arrow points from the word "input" to  $x$ .
- An arrow points from the word "parameters" to  $\theta$ .

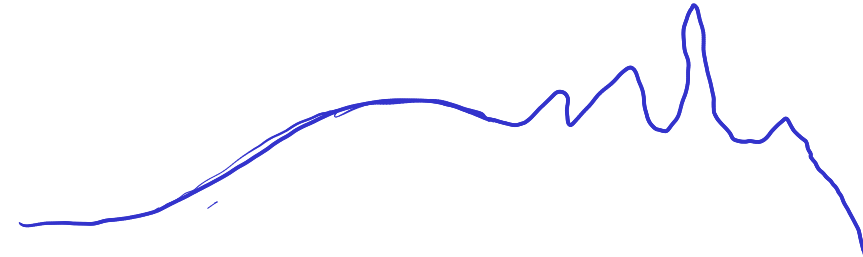
Example: Linear Function Approximation:

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

Diagram illustrating the linear function approximation equation  $f_{\theta}(x) = \theta^{\top} \beta(x)$  with handwritten blue annotations:

- An arrow points from the word "weights" to  $\theta$ .
- An arrow points from the word "features" to  $\beta(x)$ .

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



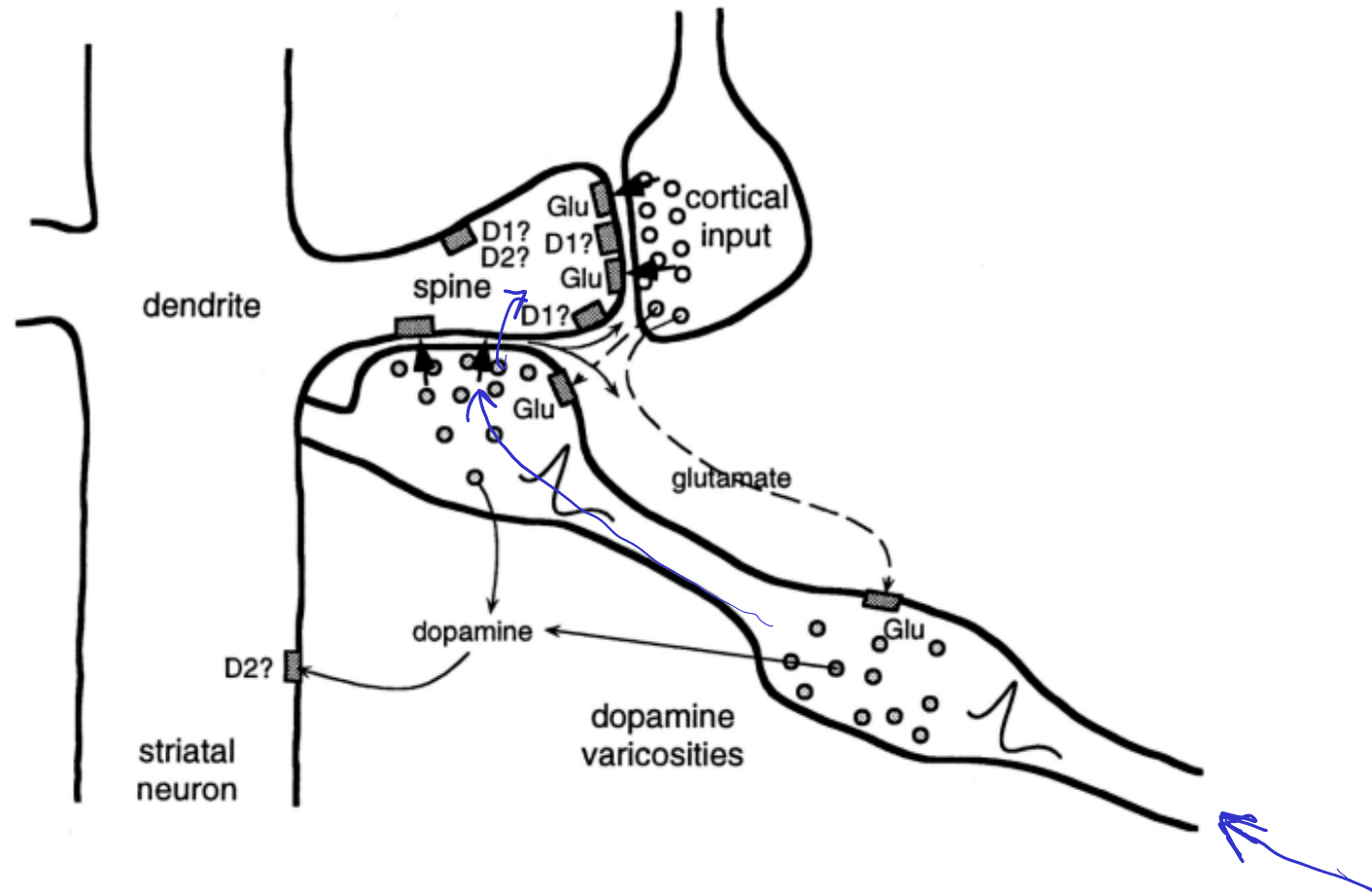
# Neural Network

# Neural Network

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

# Neural Network

$$h(x) = \sigma(\underline{W}x + b)$$



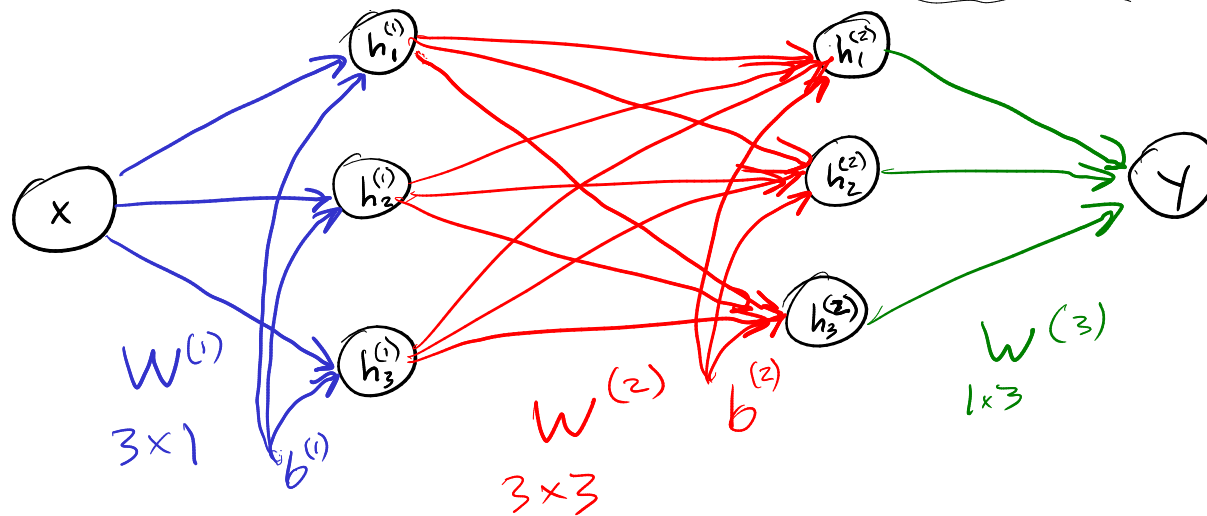
# Neural Network

# Neural Network

if  $\sigma(x) = x$

$$y = W^*x + b^*$$

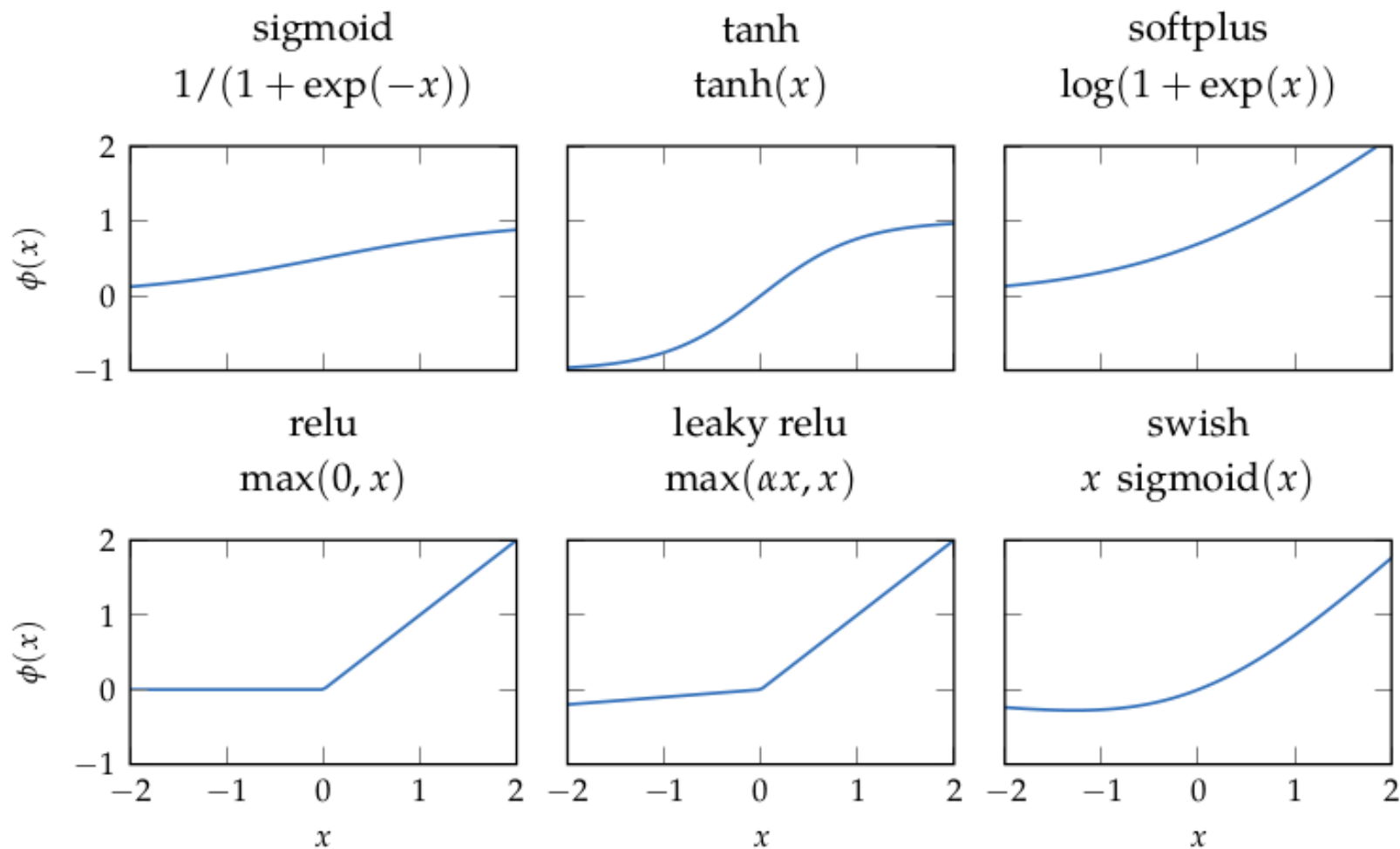
$$h(x) = \sigma(Wx + b)$$



$$f_{\theta}(x) = h^{(3)}\left(h^{(2)}\left(h^{(1)}(x)\right)\right) = W^{(3)} \sigma^{(2)}\left(W^{(2)} \sigma^{(1)}\left(W^{(1)}x + b^{(1)}\right) + b^{(2)}\right)$$

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

# Nonlinearities

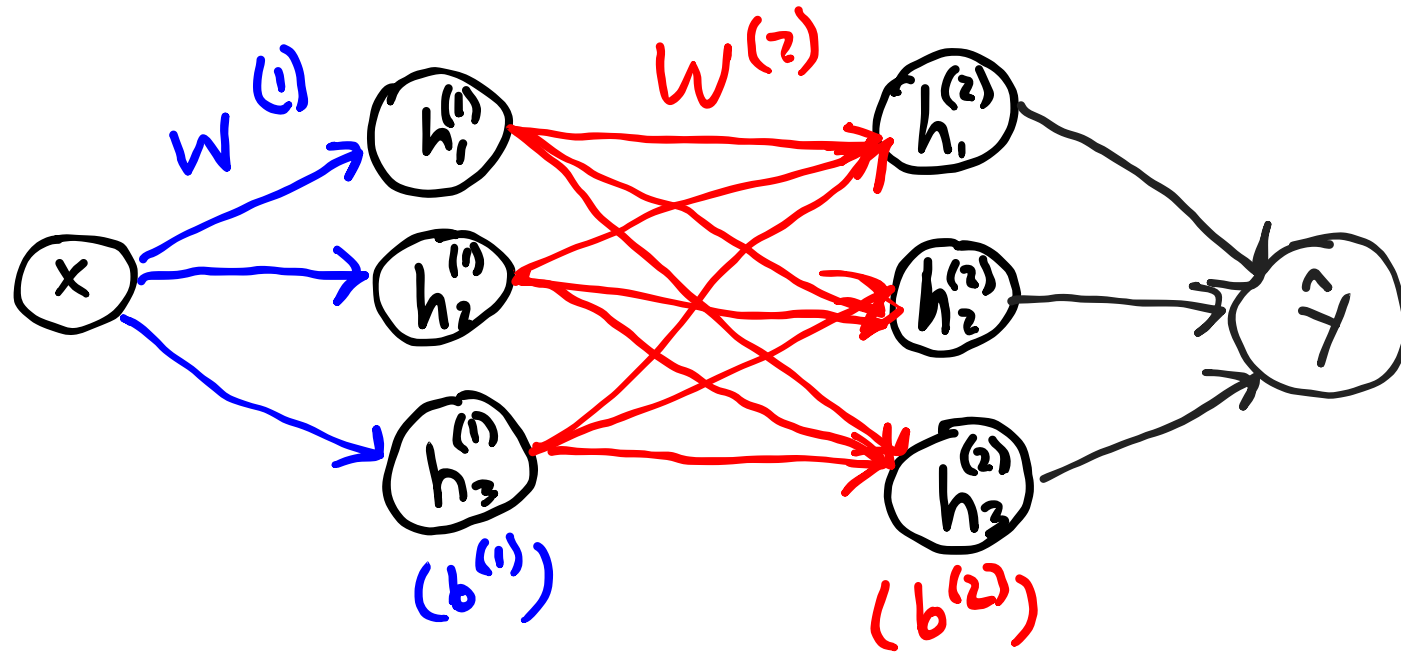




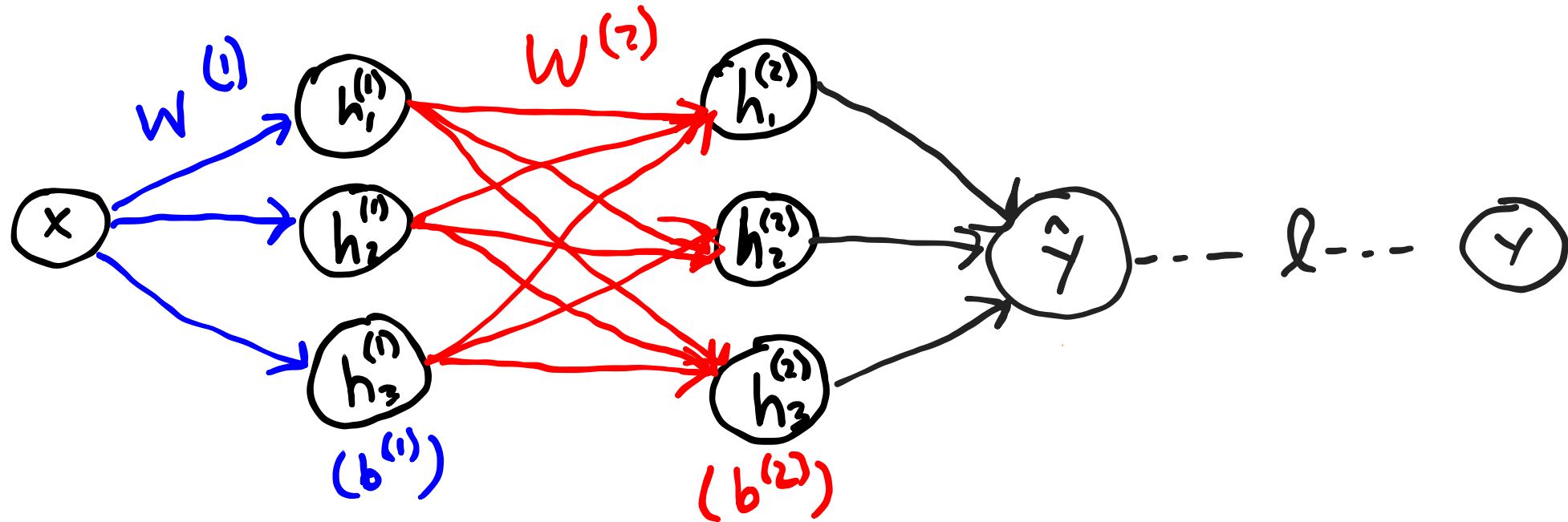
# Training

.

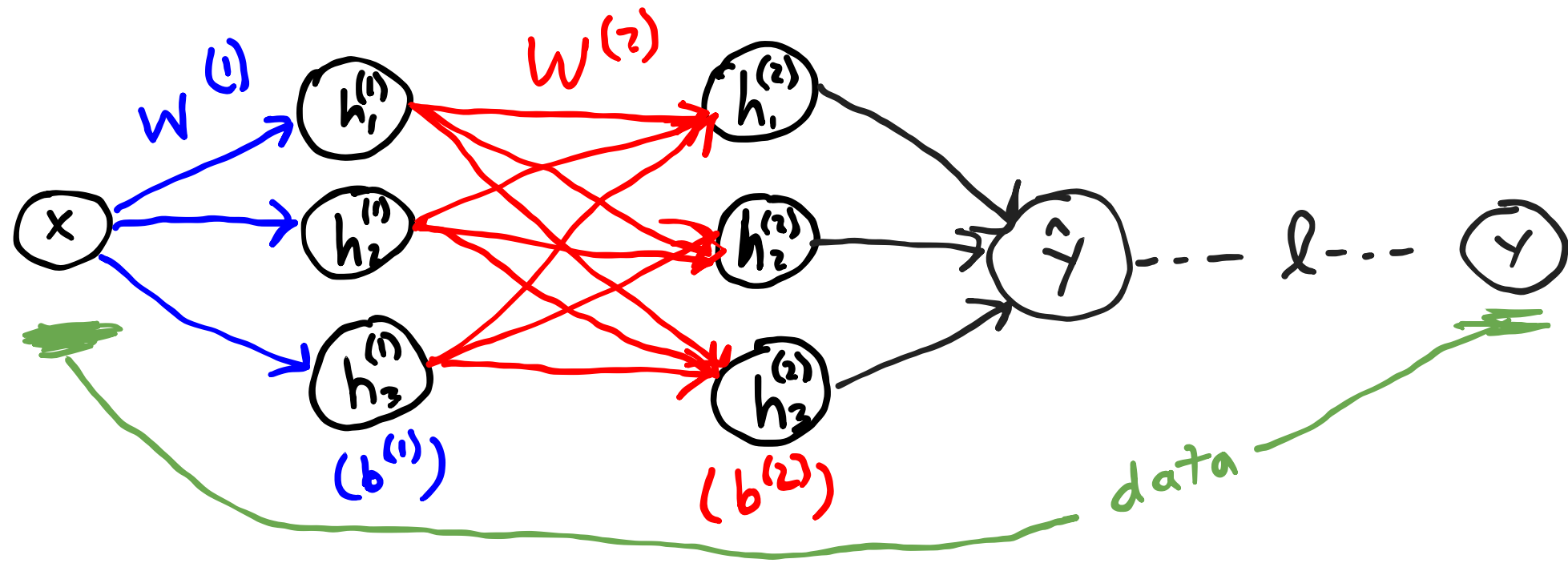
# Training



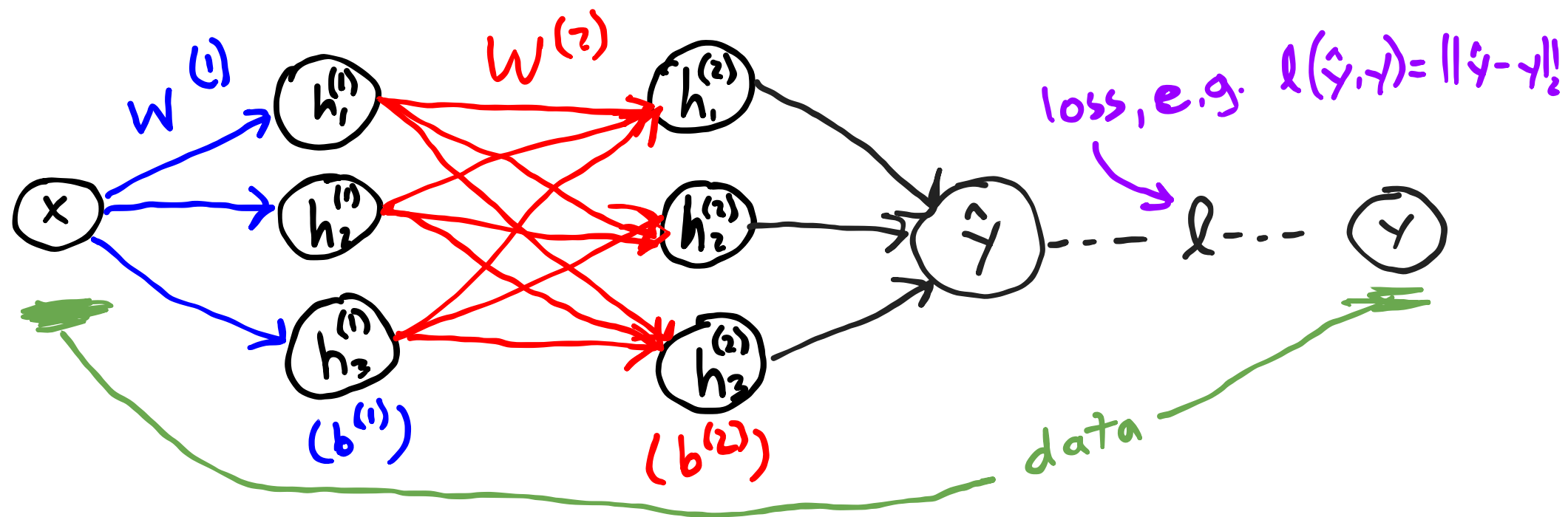
# Training



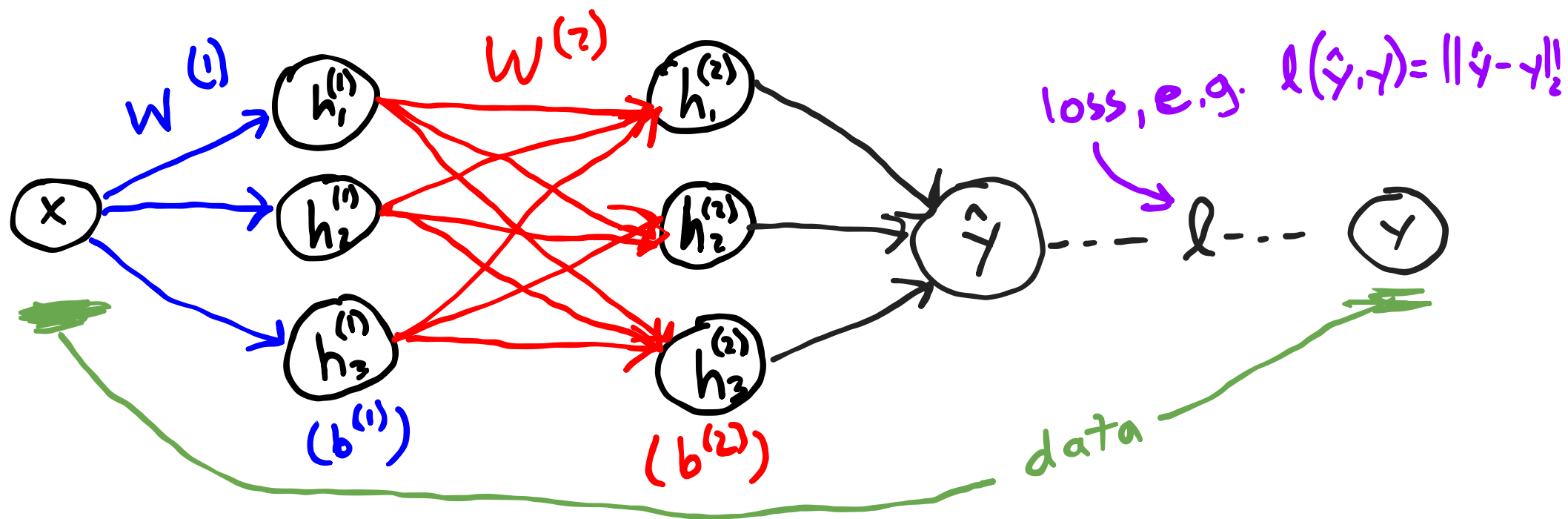
# Training



# Training

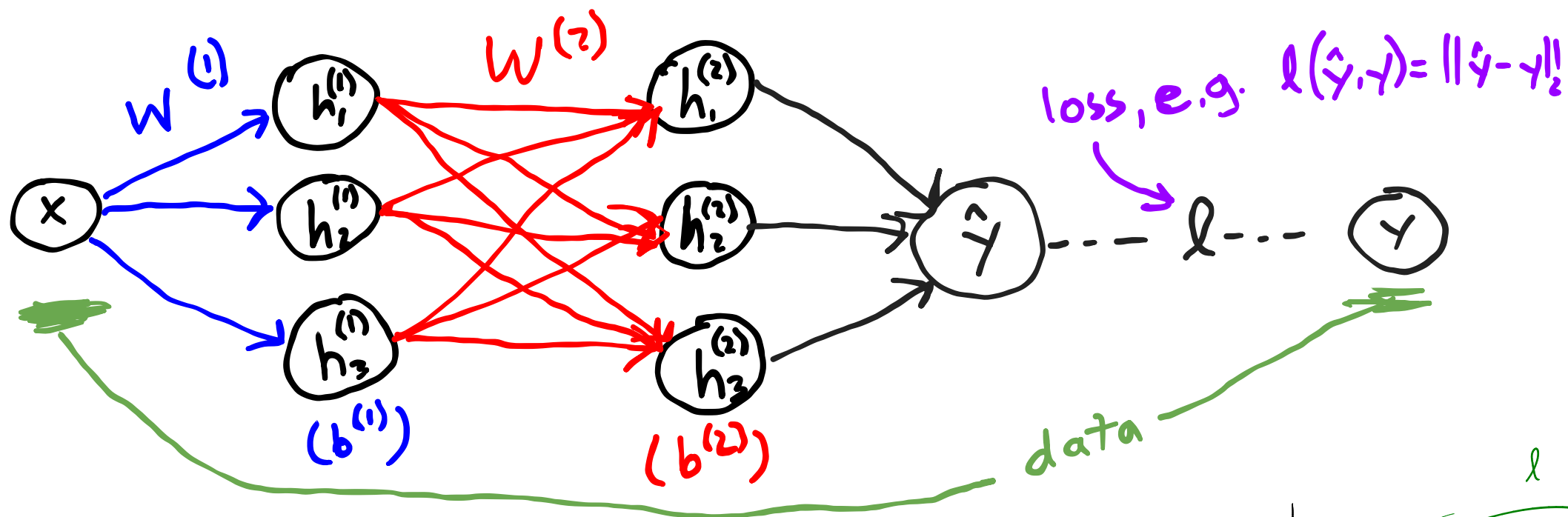


# Training



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

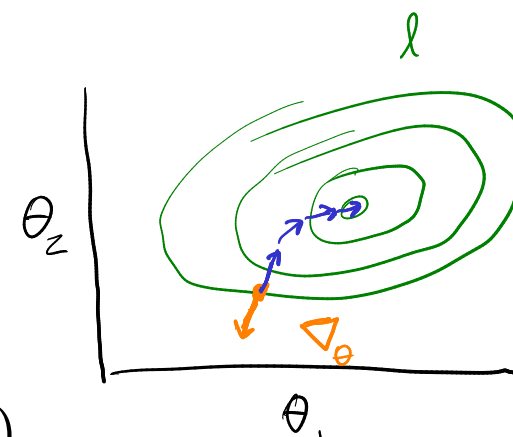
# Training



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

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

learning rate



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

$$\nabla_{\theta} l = \begin{bmatrix} \frac{\partial l}{\partial \theta_1} \\ \vdots \\ \frac{\partial l}{\partial \theta_n} \end{bmatrix}$$

# Chain Rule

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

$$l(x, y) = (f_{\theta}(x) - y)^2$$

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

$h^{(2)}$

$$\frac{\partial l}{\partial w^{(2)}} \bigg|_0 = \frac{\partial l}{\partial f} \bigg|_0 \frac{\partial f}{\partial w^{(2)}} \bigg|_0 = 2(f_{\theta}(x_0) - y_0) \cdot \underbrace{\sigma(w^{(1)} x_0 + b^{(1)})}_{\text{output of } h^{(1)}(x_0)}$$



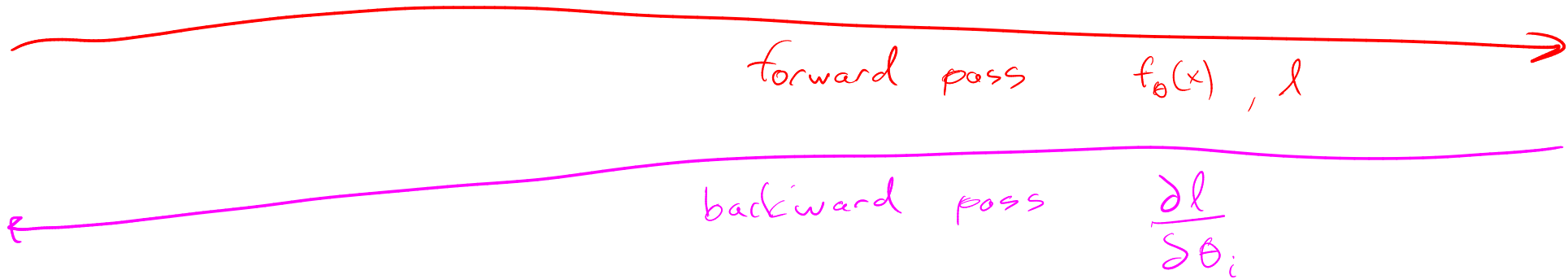
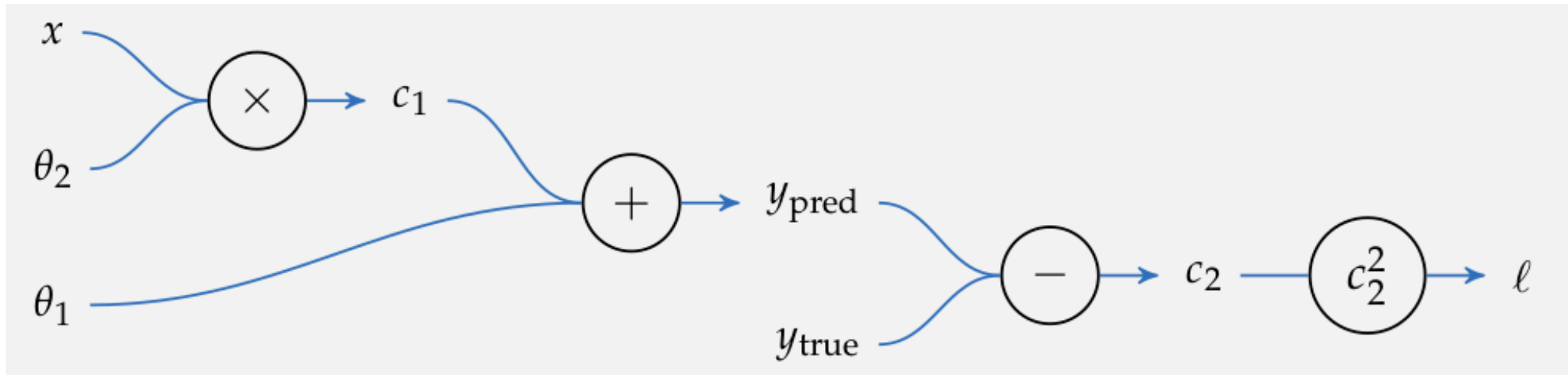
# Backprop

$$l(x, y_{\text{true}}) = (\theta_2 x + \theta_1 - y_{\text{true}})^2$$

.

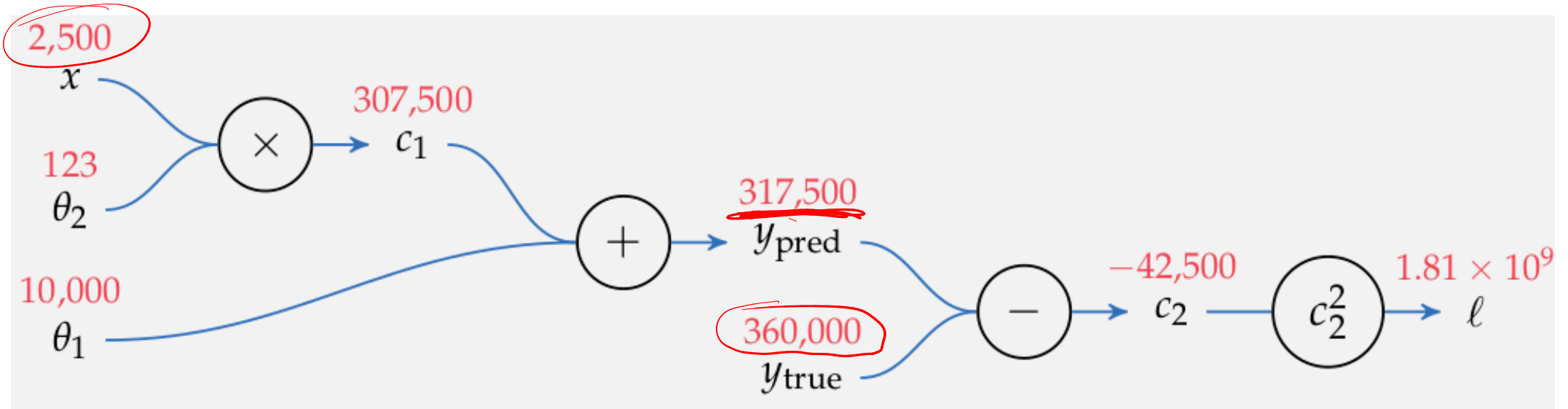
# Backprop

$$l(x, y_{\text{true}}) = (\theta_2 x + \theta_1 - y_{\text{true}})^2$$



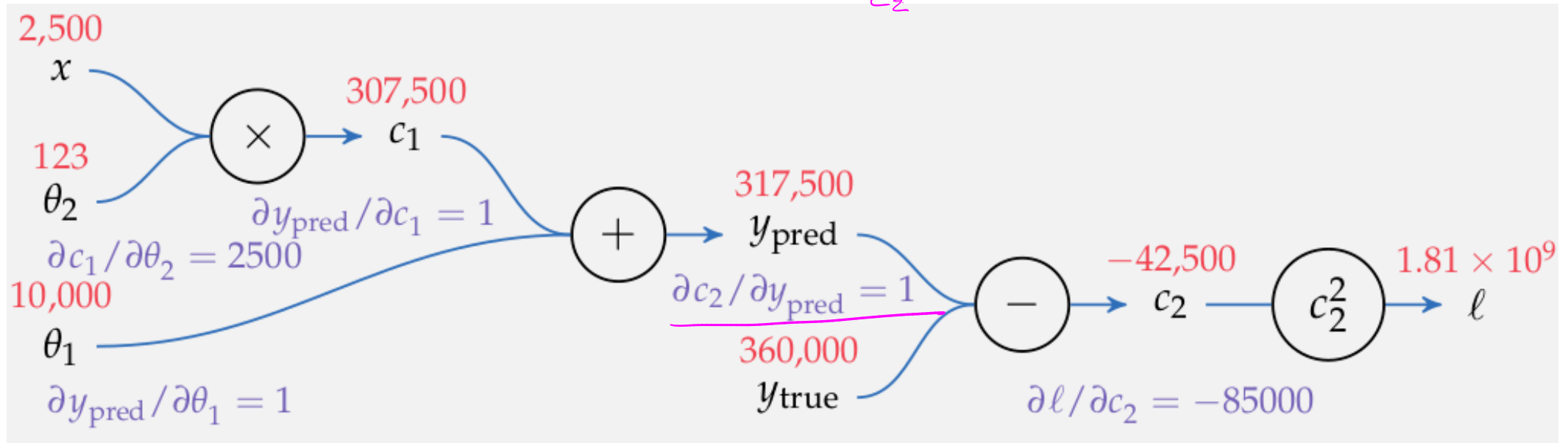
# Backprop

$$l(x, y_{\text{true}}) = (\theta_2 x + \theta_1 - y_{\text{true}})^2$$



# Backprop

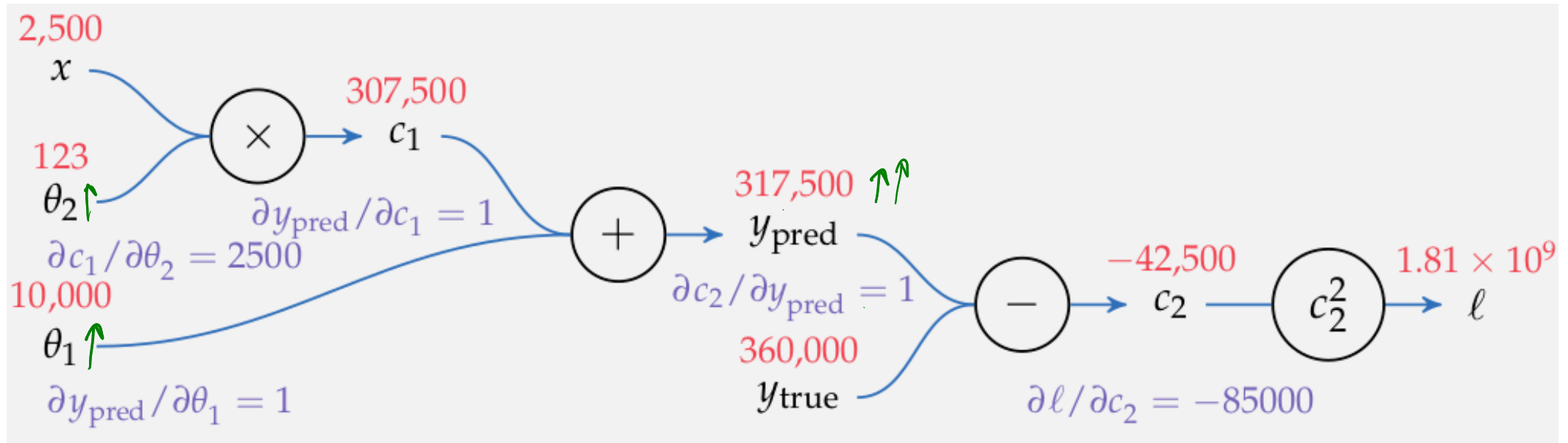
$$l(x, y_{\text{true}}) = (\underbrace{\theta_2 x + \theta_1}_{c_2} - y_{\text{true}})^2$$



$$\frac{\partial \ell}{\partial c_2} = 2 \cdot c_2 =$$

# Backprop

$$l(x, y_{\text{true}}) = (\theta_2 x + \theta_1 - y_{\text{true}})^2$$



$$\nabla_{\theta} \ell = \begin{bmatrix} \frac{\partial \ell}{\partial \theta_1} \\ \frac{\partial \ell}{\partial \theta_2} \end{bmatrix} = \begin{bmatrix} -85,000 \\ -2.125 \times 10^8 \end{bmatrix}$$

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

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

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \ell$$

```

function train(x_data, y_data;
    learning_rate=1e-3, n_epochs=1_000, save_every=50, minibatch_size=1
)
    model = Chain(
        Dense(1=>32, tanh),
        Dense(32=>32, tanh),
        Dense(32=>1)
    )
    # opt_state = Flux.setup(Adam(learning_rate), model)
    opt_state = Flux.setup(Descent(learning_rate), model)

    losses = Float32[]
    models = [deepcopy(model)]

    n_minibatches = length(y_data) ÷ minibatch_size

    @progress for epoch in 1:n_epochs
        batch_loss = zero(Float32)

        for i in 1:n_minibatches
            idxs = (1:minibatch_size) .+ minibatch_size * (i - 1)
            x_minibatch = x_data[:, idxs]
            y_minibatch = y_data[:, idxs]

            function minibatch_objective(model)
                return loss(model, x_minibatch, y_minibatch)
            end
             $l(x_0, y_0) \quad \nabla_{\theta} l$ 
            minibatch_loss, grads = Flux.withgradient(minibatch_objective, model)

            Flux.update!(opt_state, model, grads[1])
             $\theta \leftarrow \theta - \alpha \nabla_{\theta} l$ 

            batch_loss += minibatch_loss / n_minibatches
        end

        push!(losses, batch_loss)

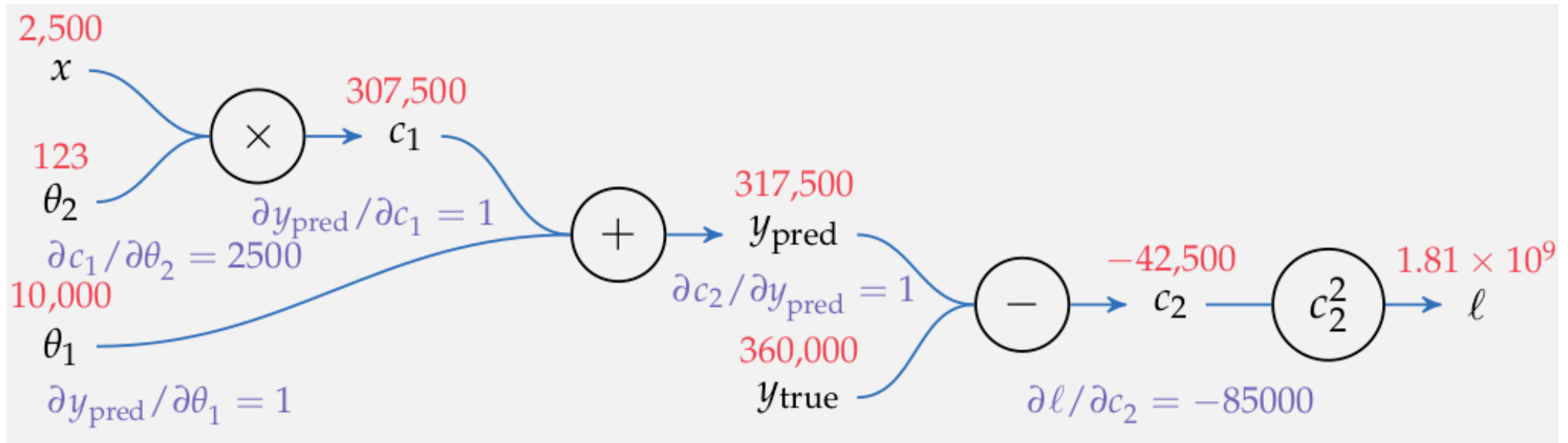
        if epoch % save_every == 0
            push!(models, deepcopy(model))
        end
    end

    return models, losses
end

```

# Backprop

$$l(x, y_{\text{true}}) = (\theta_2 x + \theta_1 - y_{\text{true}})^2$$



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

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

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

# Adaptive Step Size: ADAM

(Adaptive Moment Estimation)

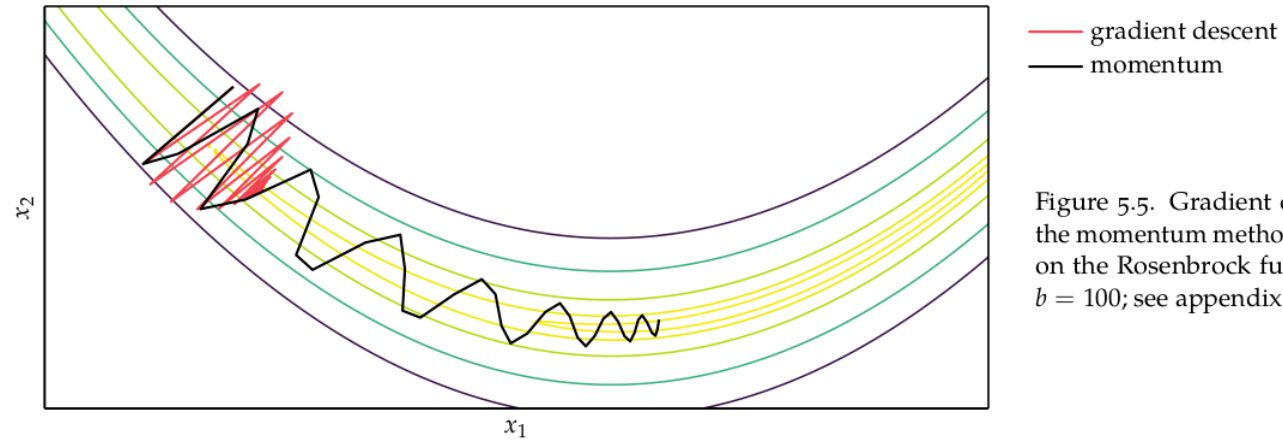


Figure 5.5. Gradient descent and the momentum method compared on the Rosenbrock function with  $b = 100$ ; see appendix B.6.

# Adaptive Step Size: ADAM

(Adaptive Moment Estimation)

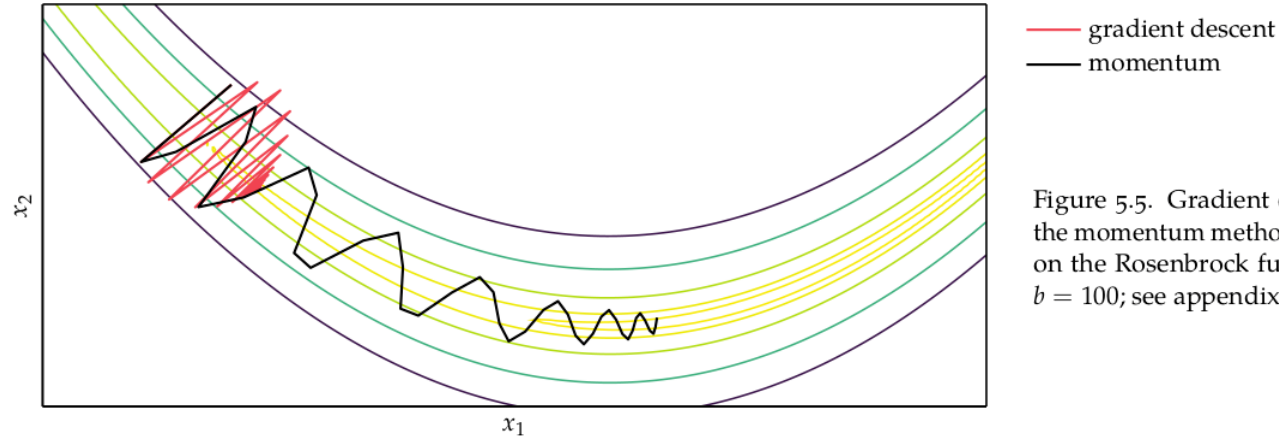


Figure 5.5. Gradient descent and the momentum method compared on the Rosenbrock function with  $b = 100$ ; see appendix B.6.

$$\text{biased decaying momentum: } \mathbf{v}^{(k+1)} = \gamma_v \mathbf{v}^{(k)} + (1 - \gamma_v) \mathbf{g}^{(k)} \quad (5.29)$$

$$\text{biased decaying sq. gradient: } \mathbf{s}^{(k+1)} = \gamma_s \mathbf{s}^{(k)} + (1 - \gamma_s) (\mathbf{g}^{(k)} \odot \mathbf{g}^{(k)}) \quad (5.30)$$

$$\text{corrected decaying momentum: } \hat{\mathbf{v}}^{(k+1)} = \mathbf{v}^{(k+1)} / (1 - \gamma_v^k) \quad (5.31)$$

$$\text{corrected decaying sq. gradient: } \hat{\mathbf{s}}^{(k+1)} = \mathbf{s}^{(k+1)} / (1 - \gamma_s^k) \quad (5.32)$$

$$\text{next iterate: } \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha \hat{\mathbf{v}}^{(k+1)} / \left( \epsilon + \sqrt{\hat{\mathbf{s}}^{(k+1)}} \right) \quad (5.33)$$

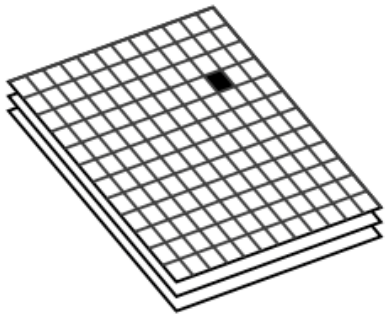
$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \ell$$

<sup>12</sup> According to the original paper, good default settings are  $\alpha = 0.001$ ,  $\gamma_v = 0.9$ ,  $\gamma_s = 0.999$ , and  $\epsilon = 1 \times 10^{-8}$ .

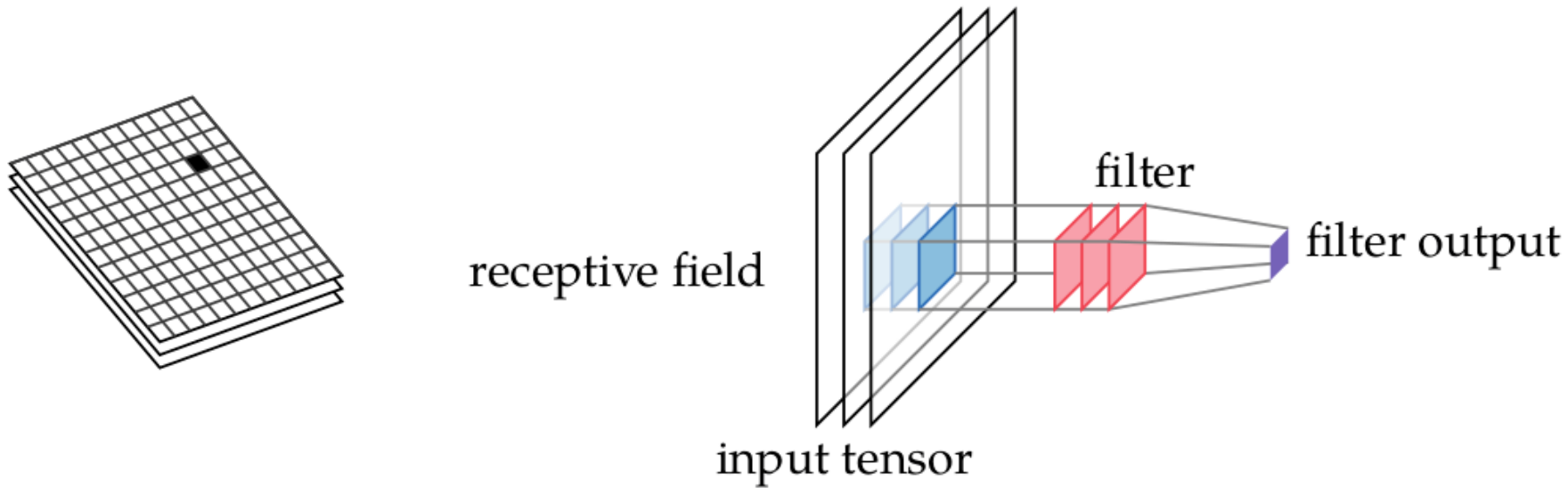
$\odot$  means elementwise multiplication.

# On Your Radar: ConvNets

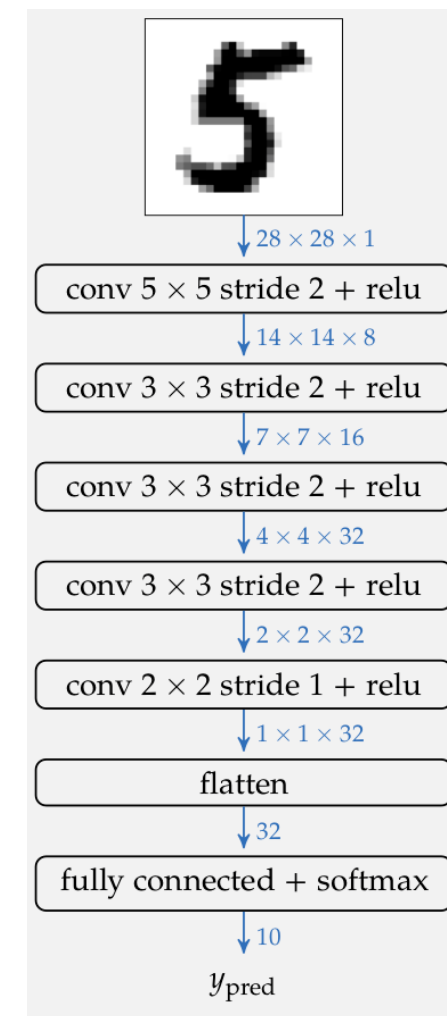
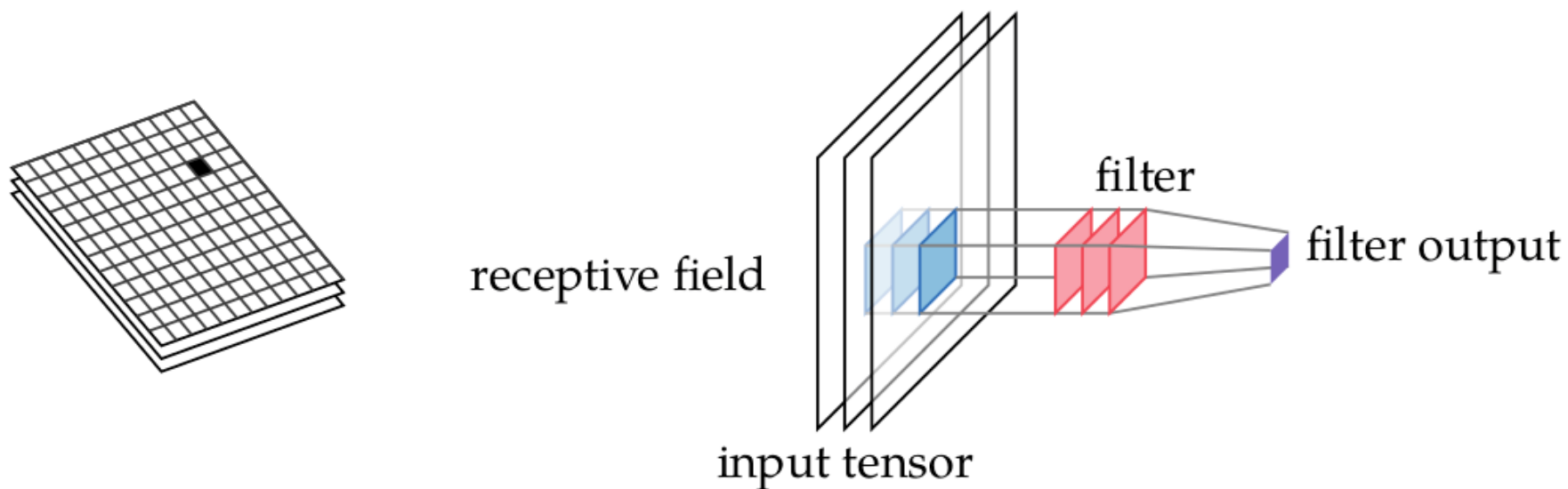
# On Your Radar: ConvNets



# 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

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

e.g. Batch norm, layer norm, dropout

# On Your Radar: Skip Connections (Resnets)

# Resources

OpenAI Spinning up