



Loss and Optimization

Loss Functions

- How “bad” is our model?

$$L(\theta) = \sum_i \ell(\theta; x_i, y_i)$$

Example: least-squares regression

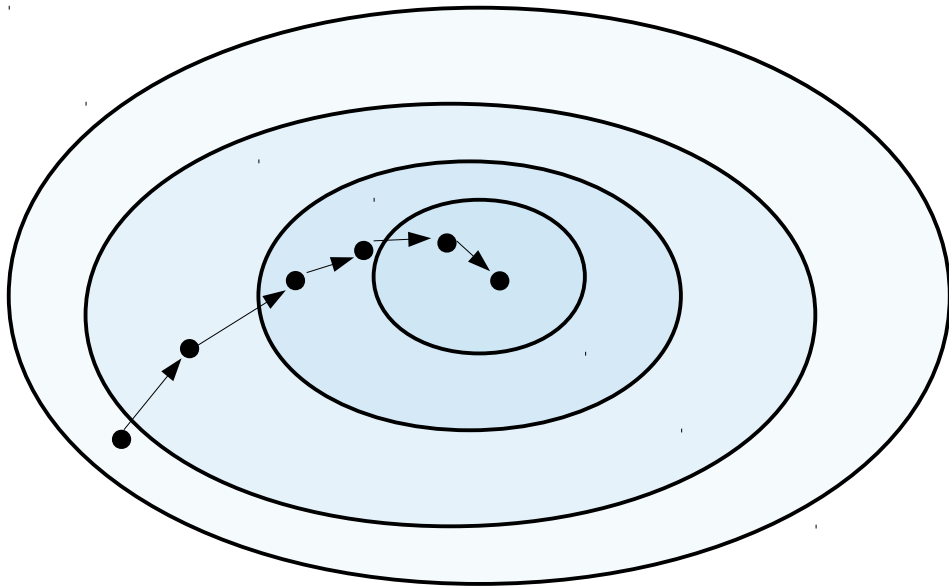
$$L(\theta) = \sum_i (y_i - f(\theta; x_i))^2$$

Classification: Negative log-likelihood

$$L(\theta) = \sum_i -\log(p(y_i | \theta, x_i))$$

Optimization – Gradient Descent

How do we find $\operatorname{argmin}_{\theta} L(\theta)$



Initialize random θ

N iterations:

Compute $\nabla_{\theta} L(\theta)$

$\theta \leftarrow \theta - \epsilon \nabla_{\theta} L(\theta)$

Gradients

- How do we compute $\nabla_{\theta} L(\theta)$?

$$\nabla_{\theta} L(\theta) = \nabla_{\theta} \left(\sum_i \ell(\theta | x_i, y_i) \right) = \sum_i \nabla_{\theta} \ell(\theta | x_i, y_i)$$

Squared error:

$$\ell(\theta | x_i, y_i) = ((\mathbf{W} \mathbf{x}_i + \mathbf{b}) - y_i)^2$$
$$\nabla_{\mathbf{W}} ((\mathbf{W} \mathbf{x}_i + \mathbf{b}) - y_i)^2 = 2((\mathbf{W} \mathbf{x}_i + \mathbf{b}) - y_i) \mathbf{x}_i$$
$$\nabla_{\mathbf{b}} ((\mathbf{W} \mathbf{x}_i + \mathbf{b}) - y_i)^2 = 2((\mathbf{W} \mathbf{x}_i + \mathbf{b}) - y_i)$$

Negative log-likelihood: $\ell(\theta | x_i, y_i) = -\log(\text{softmax}(\mathbf{W} \mathbf{x}_i + \mathbf{b}))_{y_i}$

$$\nabla_{\mathbf{W}_j} \left(-\log(\text{softmax}(\mathbf{W} \mathbf{x}_i + \mathbf{b}))_{y_i} \right) = (\text{softmax}(\mathbf{W} \mathbf{x}_i + \mathbf{b})_j - [y_i = j]) \mathbf{x}_i$$

$$\nabla_{\mathbf{b}_j} \left(-\log(\text{softmax}(\mathbf{W} \mathbf{x}_i + \mathbf{b}))_{y_i} \right) = \text{softmax}(\mathbf{W} \mathbf{x}_i + \mathbf{b})_j - [y_i = j]$$

Gradients – Computation Graphs

$$(y - (Wx + b))^2$$

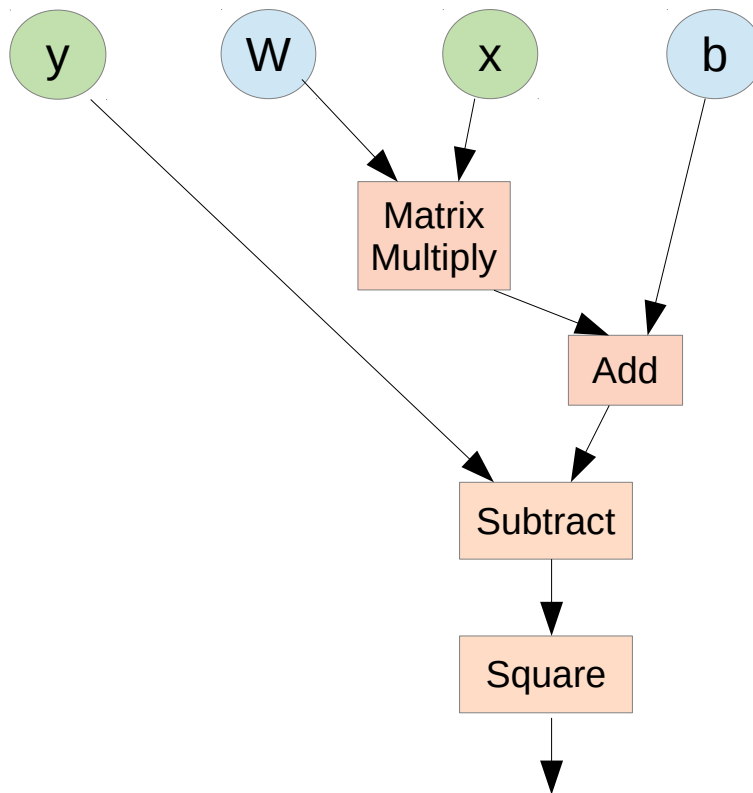
$$\text{sq}(\text{sub}(y, \text{add}(\text{matmul}(W, x), b)))$$

$$\frac{\partial}{\partial W} \text{sq}(\text{sub}(y, \text{add}(\text{matmul}(W, x), b)))$$

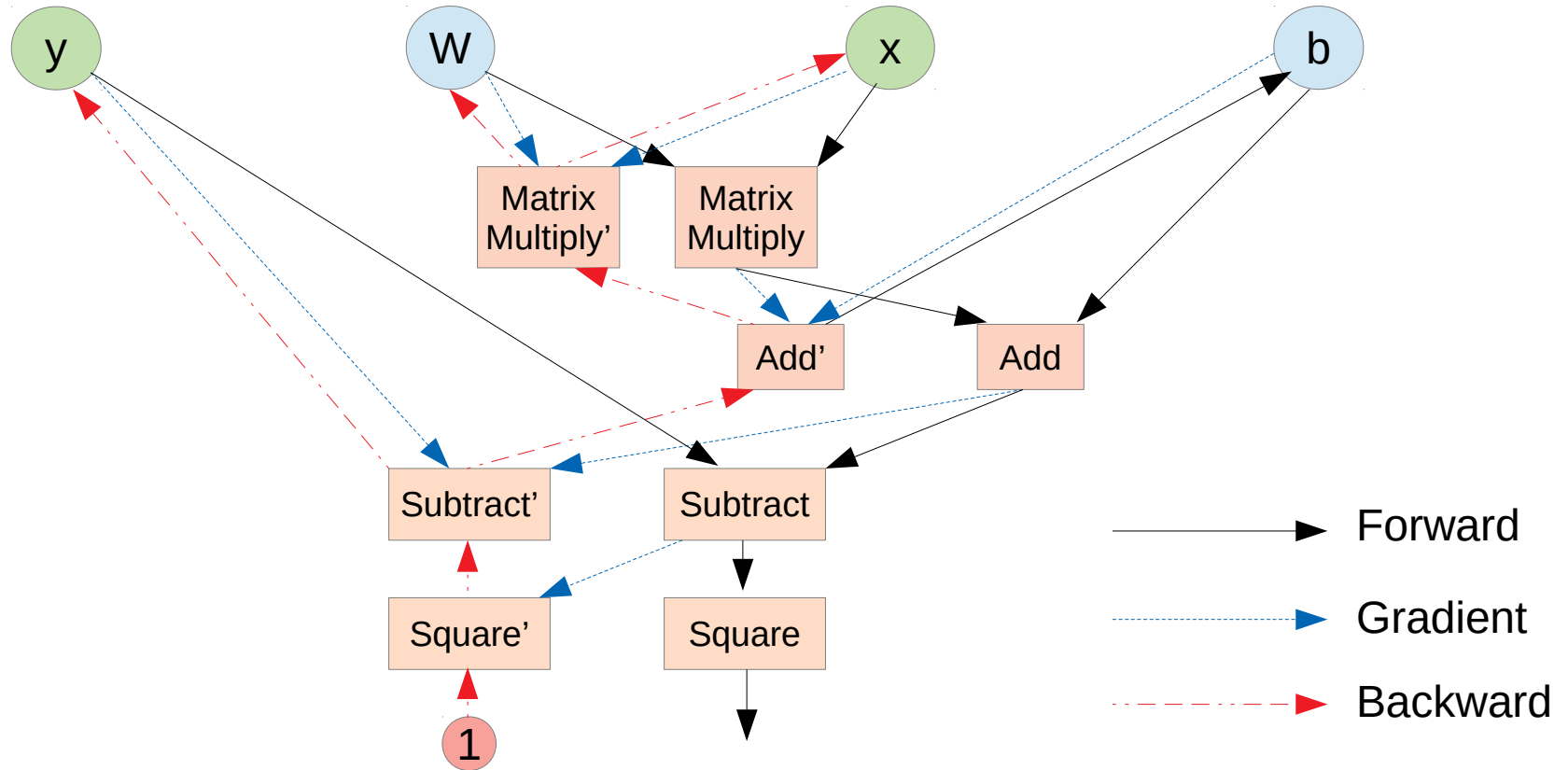
$$= \text{sq}'(\text{sub}(y, \text{add}(\text{matmul}(W, x), b))) * \frac{\partial}{\partial W} \text{sub}(y, \text{add}(\text{matmul}(W, x), b))$$

= ...

$$= \text{sq}'(\text{sub}(y, \text{add}(\text{matmul}(W, x), b))) * \text{sub}'(y, \text{add}(\text{matmul}(W, x), b)) * \text{add}'(\text{matmul}(W, x), b) * \text{matmul}'(W, x)$$

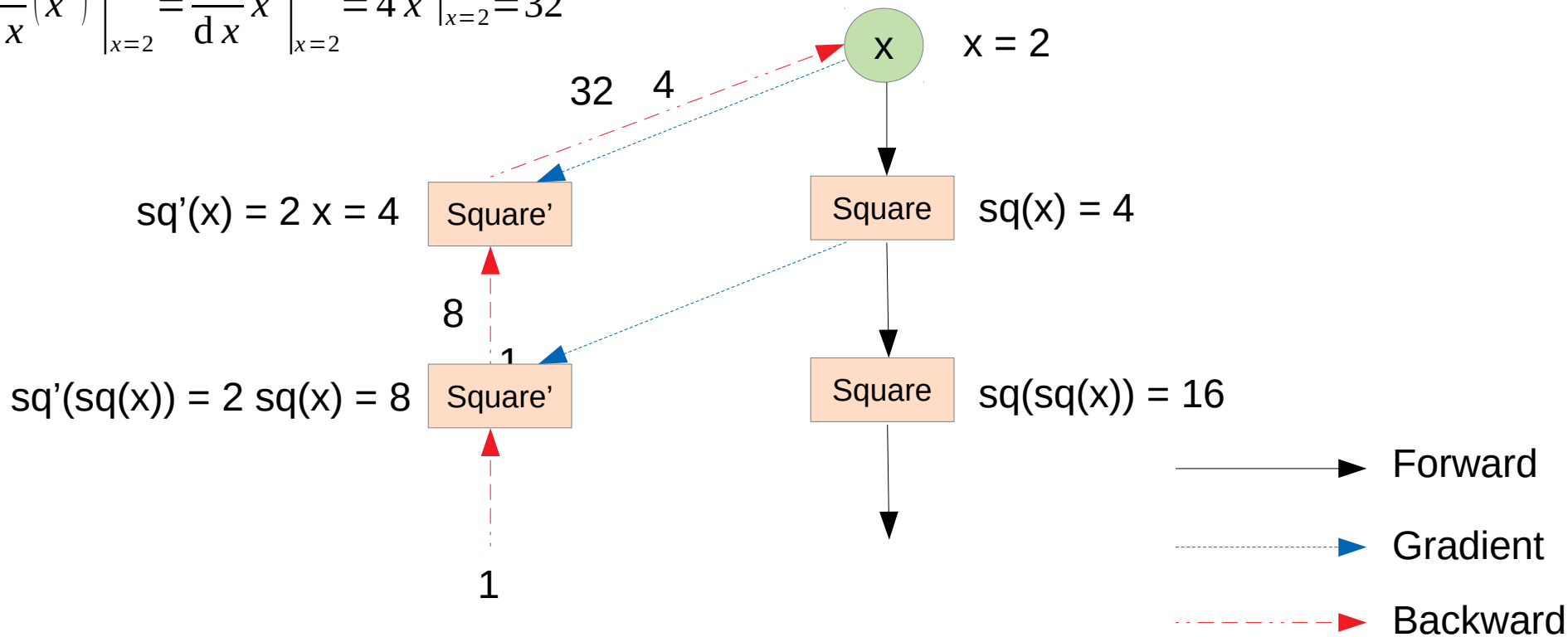


Gradients – Computation Graphs



Gradients – Computation Graphs

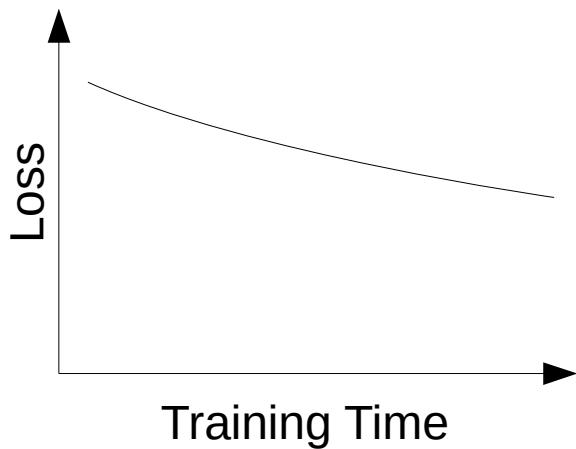
$$\left. \frac{d}{dx} (x^2)^2 \right|_{x=2} = \left. \frac{d}{dx} x^4 \right|_{x=2} = 4x^3 \Big|_{x=2} = 32$$



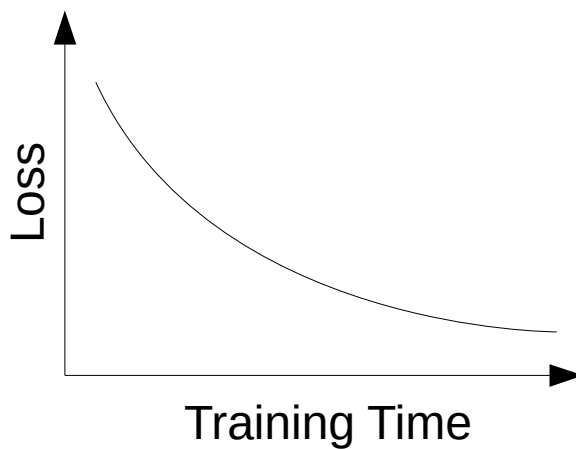
Learning Rate

$$\theta \leftarrow \theta - \epsilon \nabla_{\theta} L(\theta)$$

Too small



Ideal



Too large

