

Gradient Descent

Numerical derivative

Random Search

Save weight if current loss better than best loss

Based on MS DL Course of Boris Zubarev @bobazooba

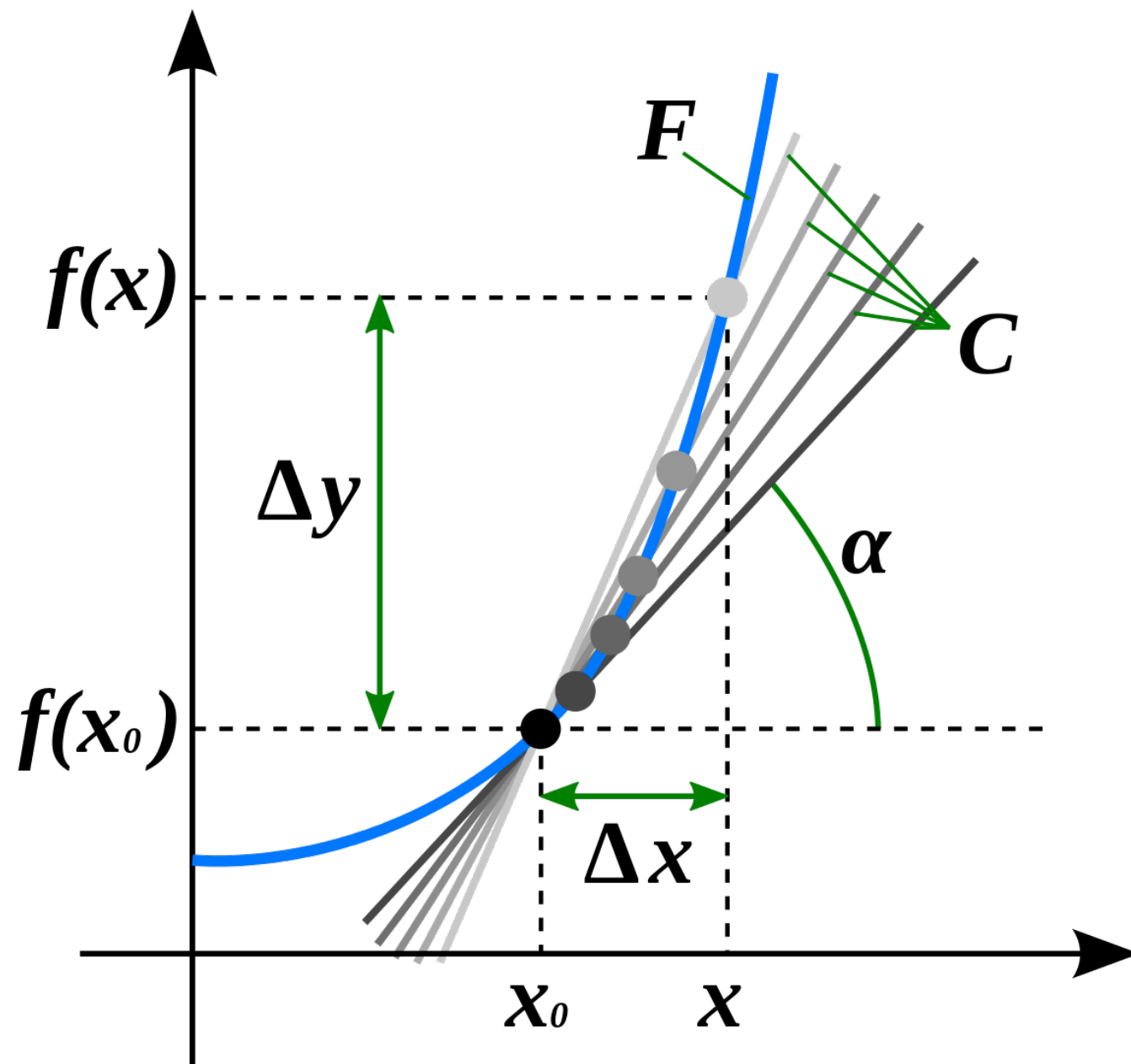
Gradient Descent

Numerical derivative

$$\frac{\partial L}{\partial w} \approx \frac{L(w + \epsilon) - L(w - \epsilon)}{2\epsilon}$$

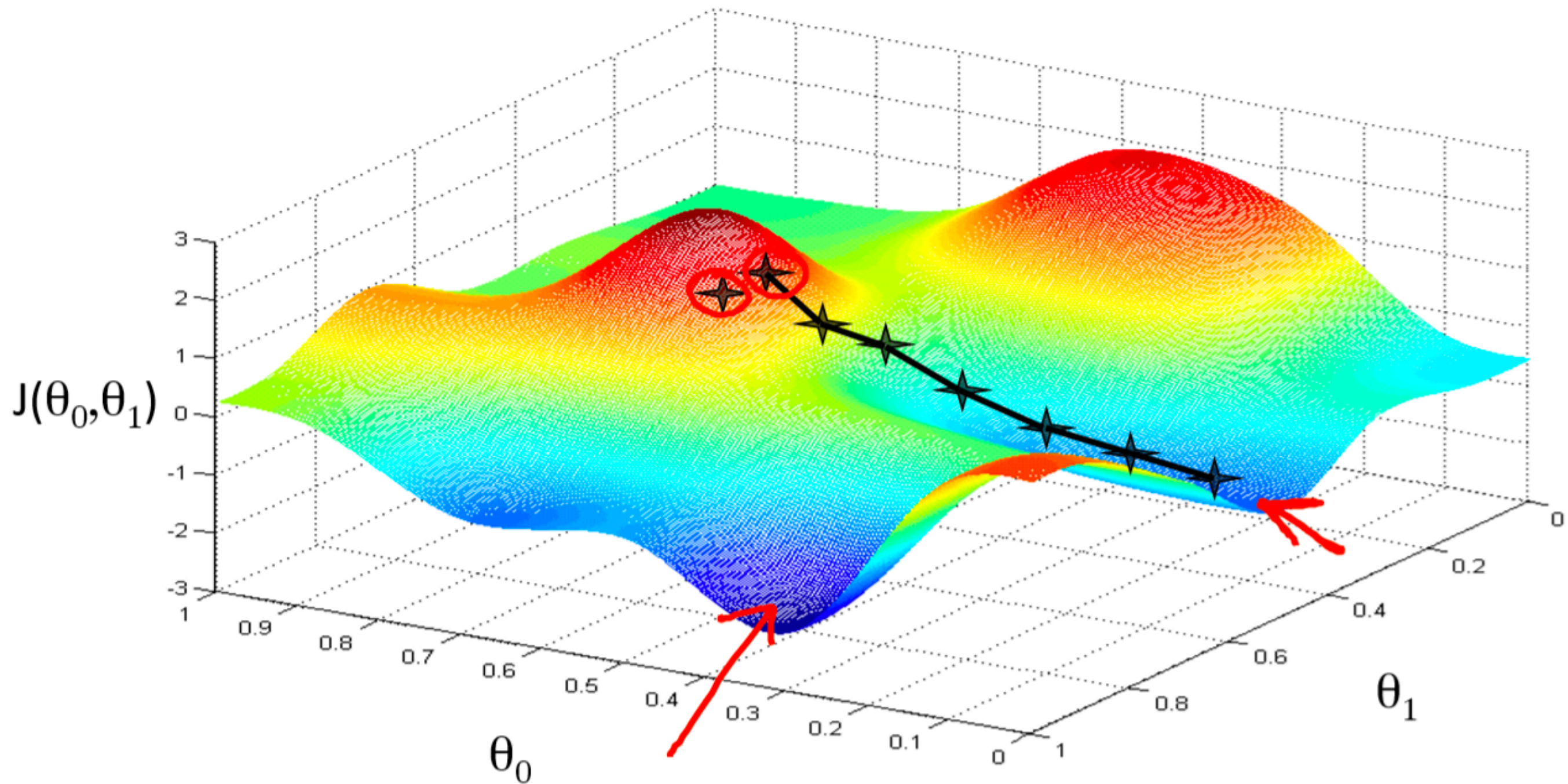
Gradient Descent

Derivative



Gradient Descent

Local minima

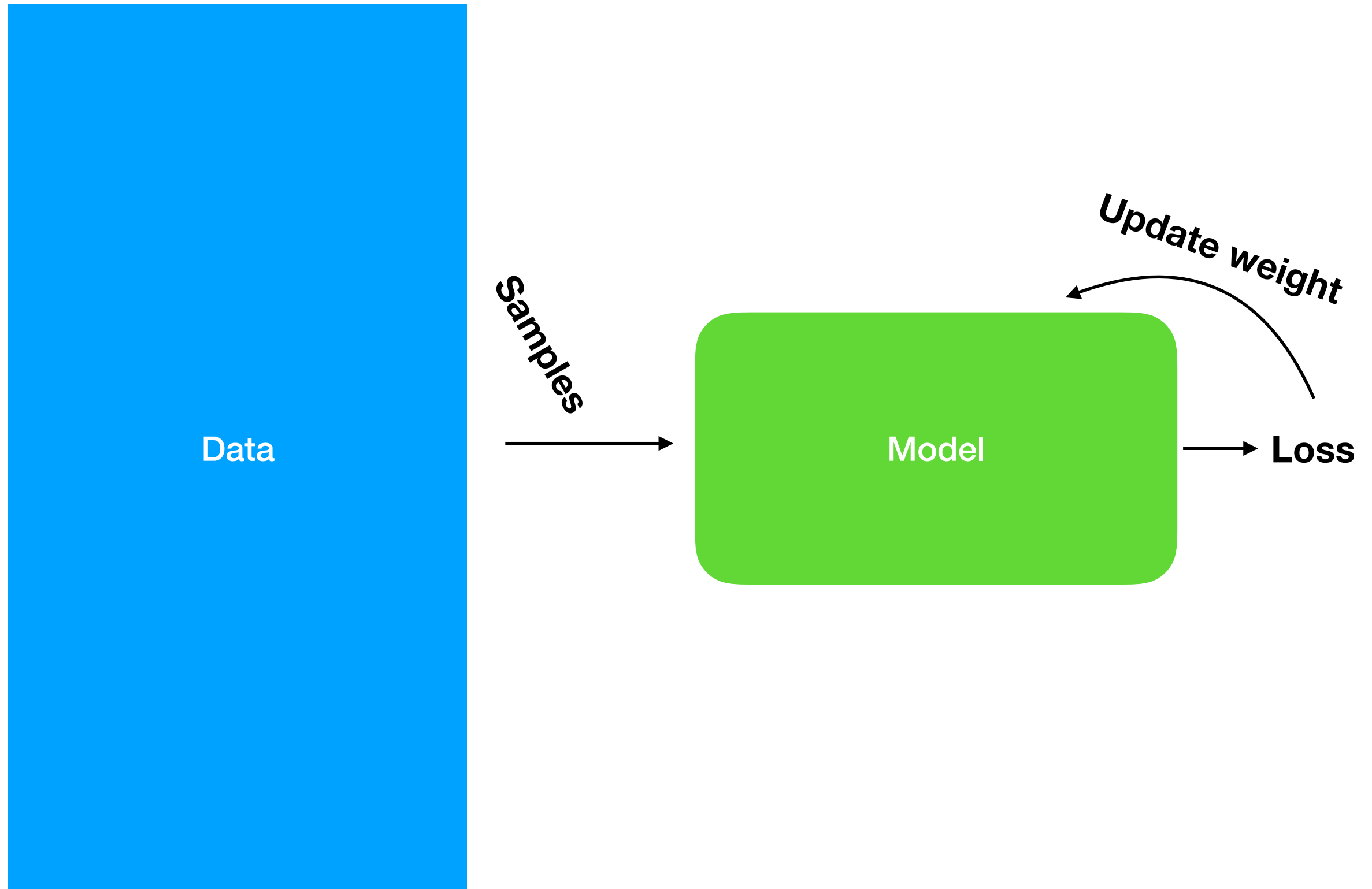


Gradient Descent

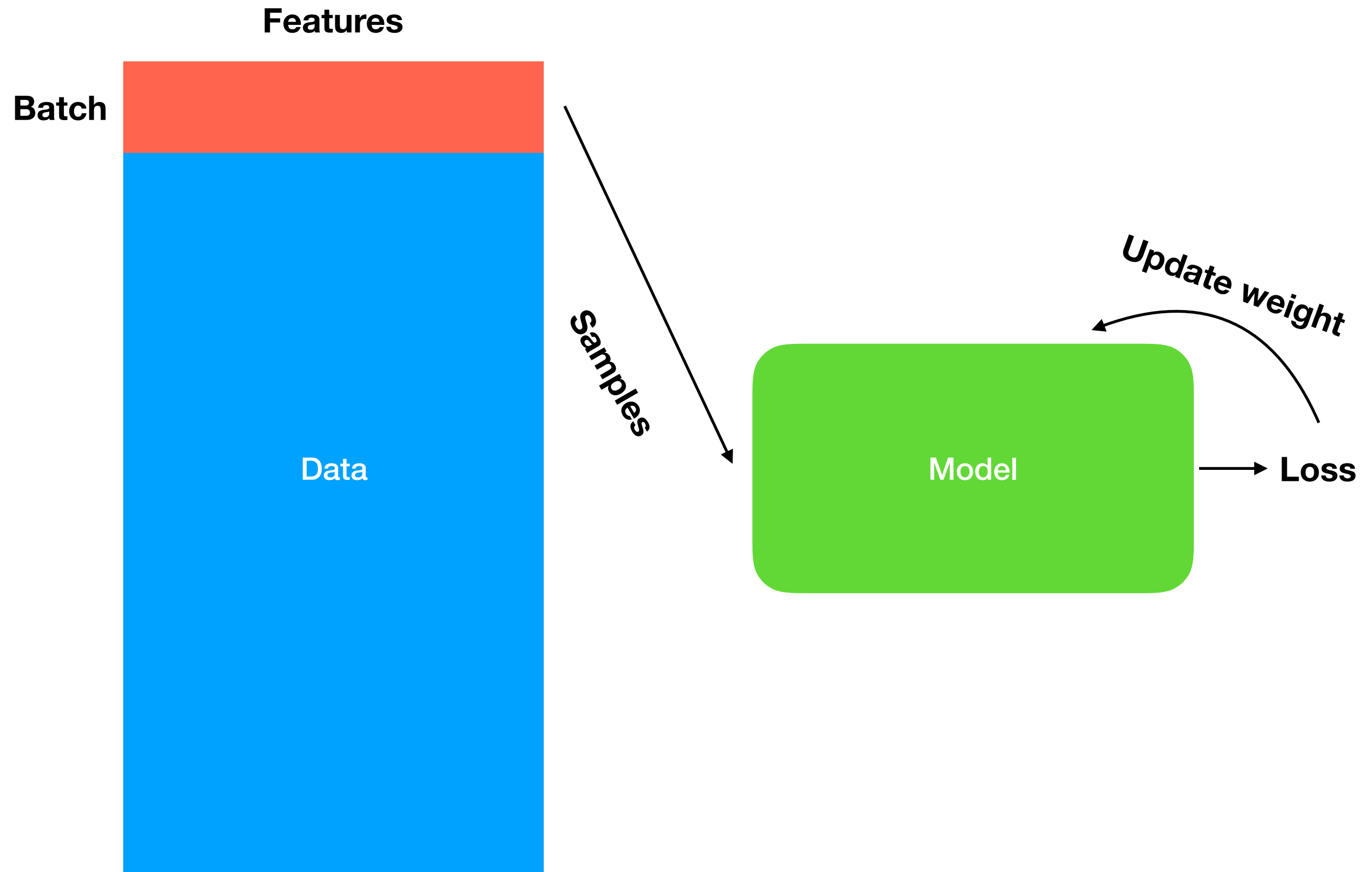
```
W = W - learning_rate * dLdW  
b = b - learning_rate * dLdb
```

Gradient Descent

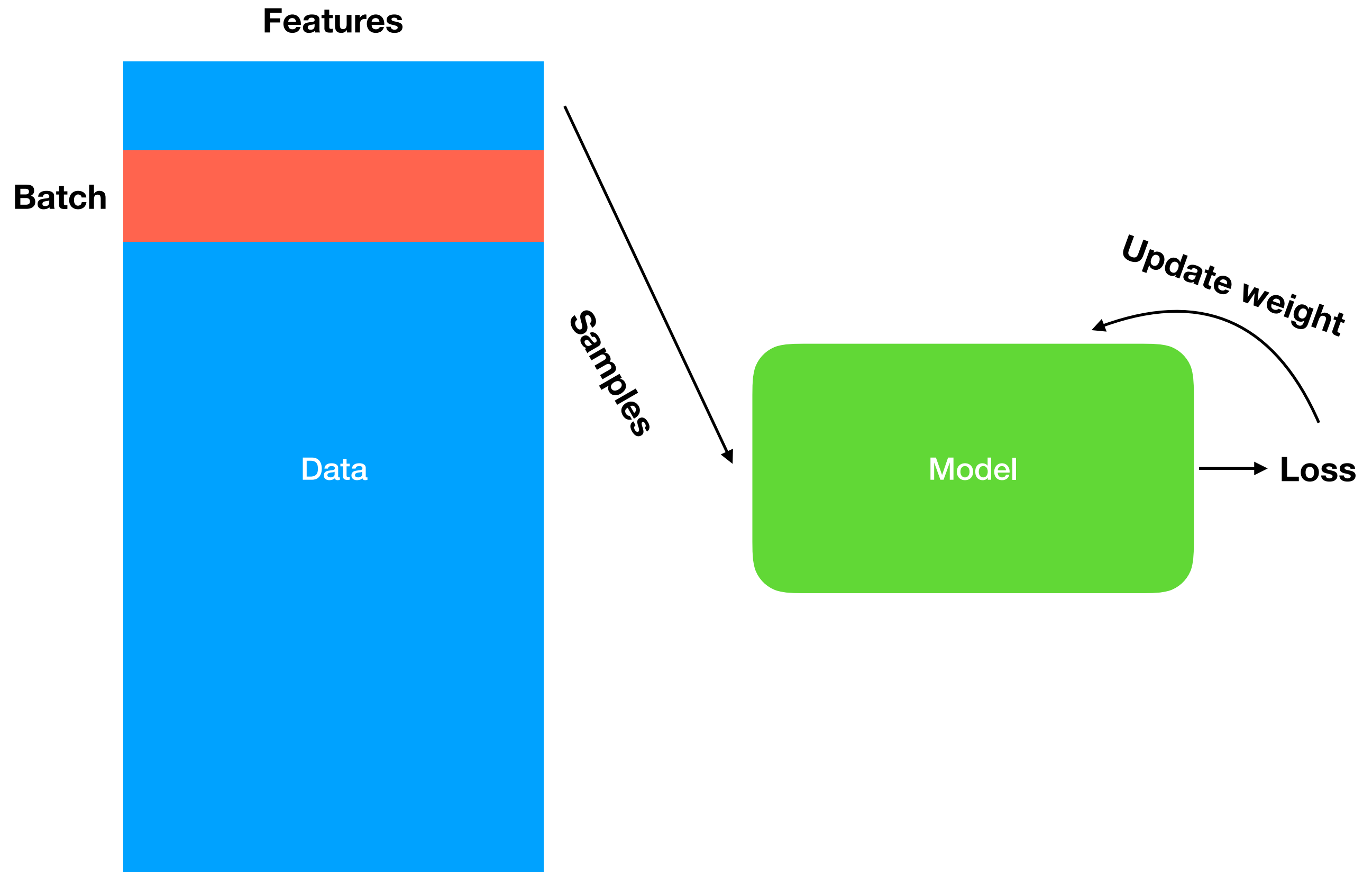
Features



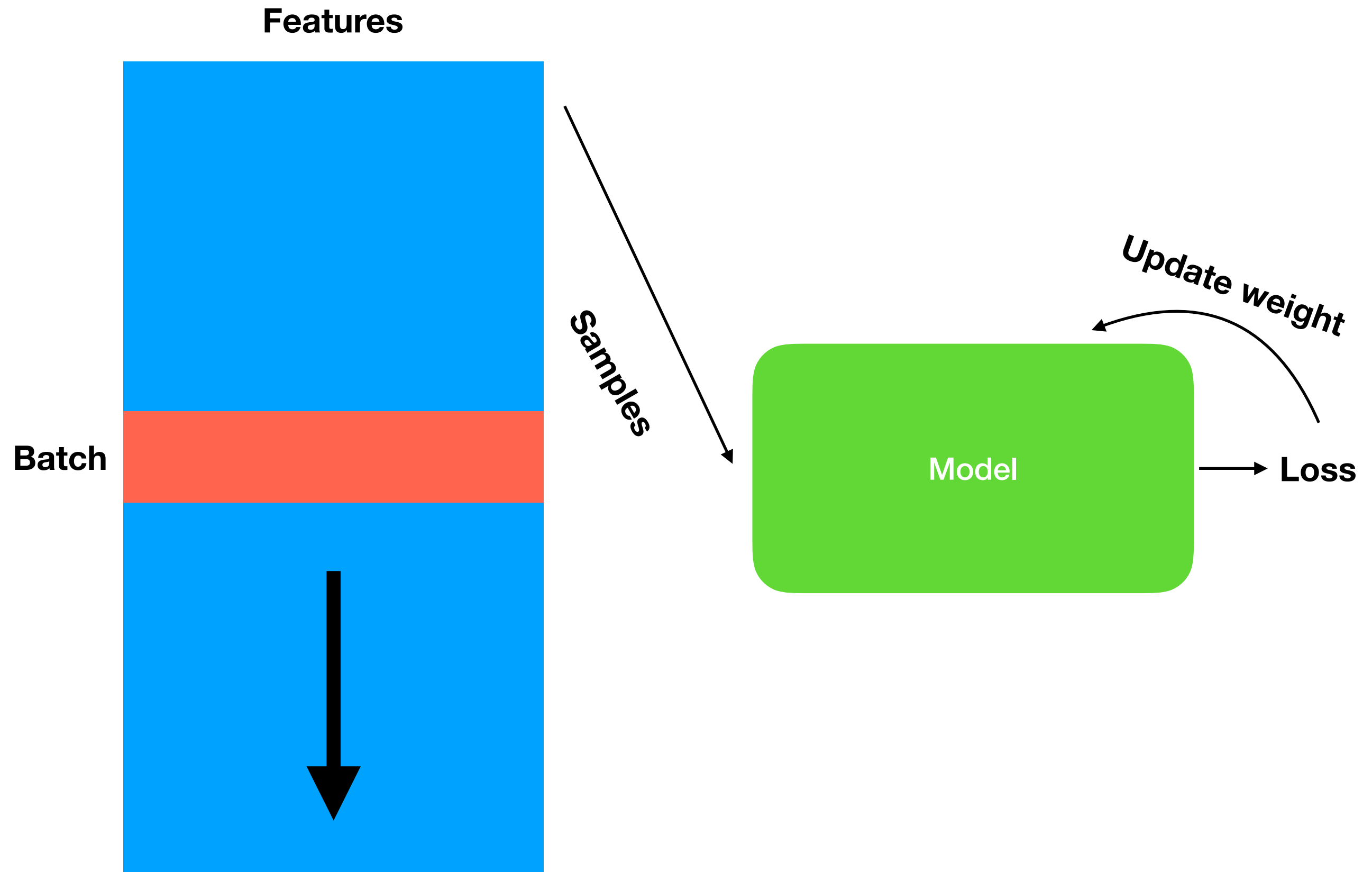
Stochastic Gradient Descent



Stochastic Gradient Descent

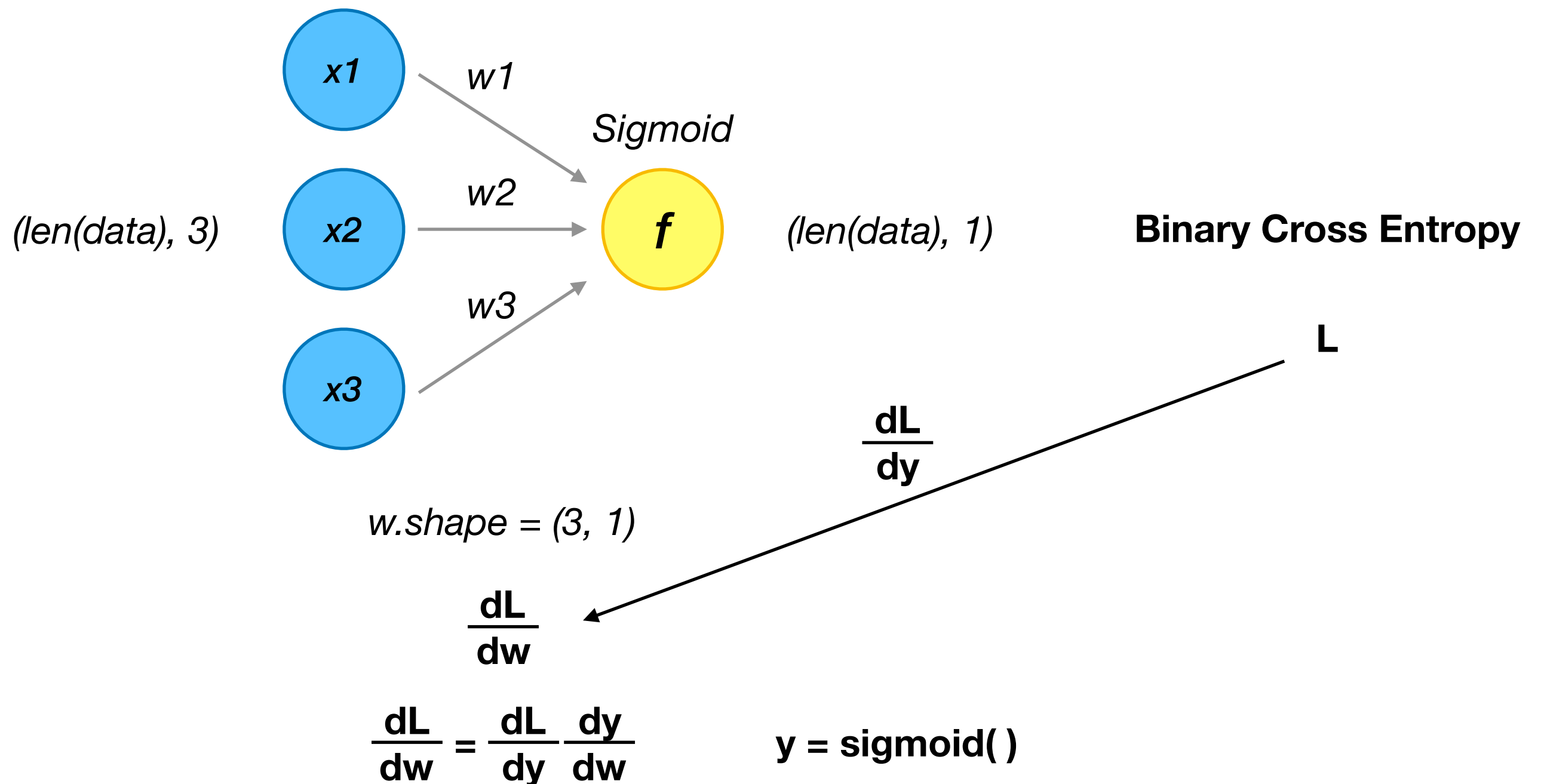


Stochastic Gradient Descent



Logistic Regression

Train

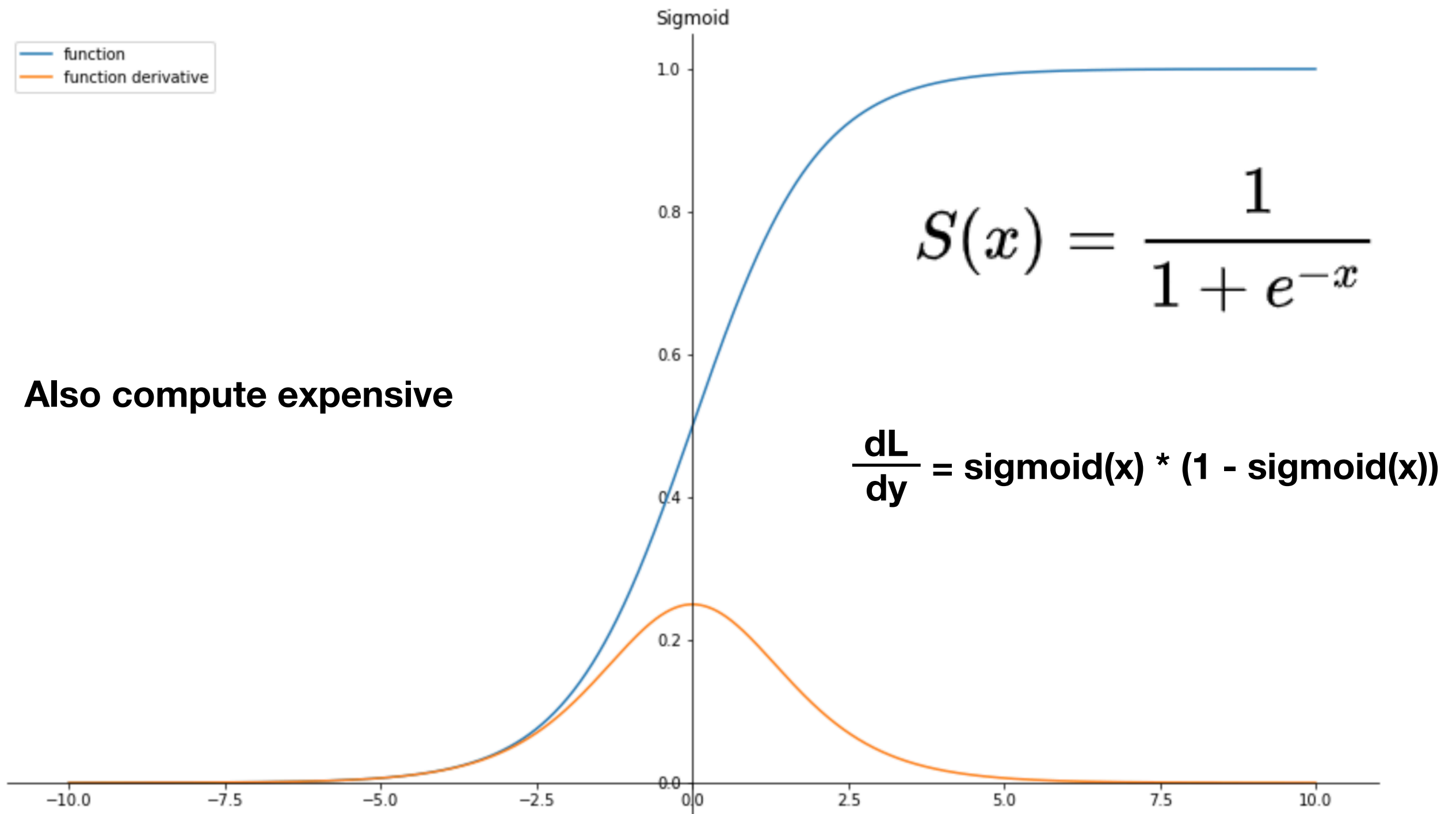


Regularization

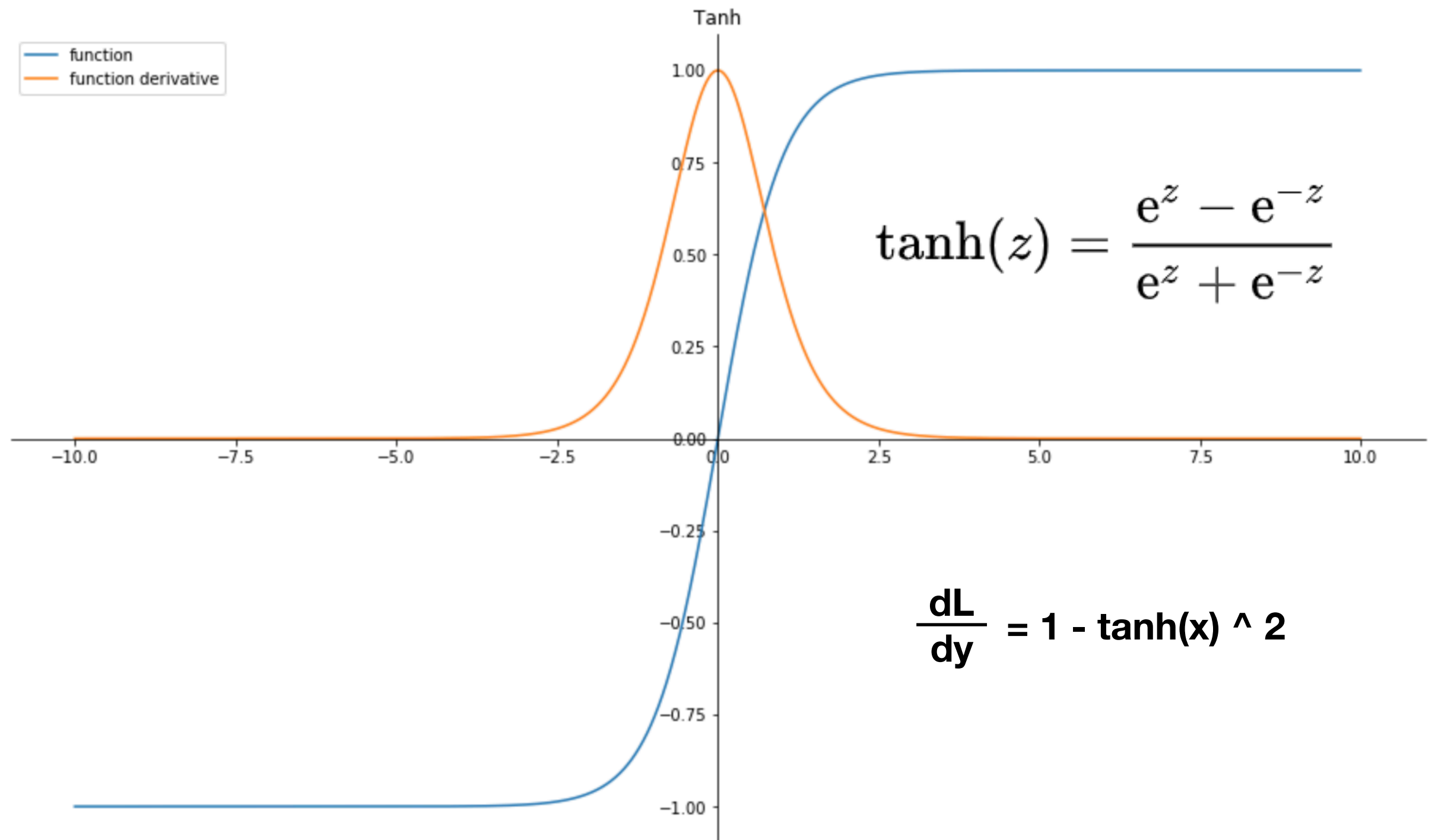
$$\text{BCELoss} = - (y * \log(\text{pred}) + (1 - y) * \log(1 - \text{pred})) + \text{reg} * R(W)$$

reg - hyperparameter

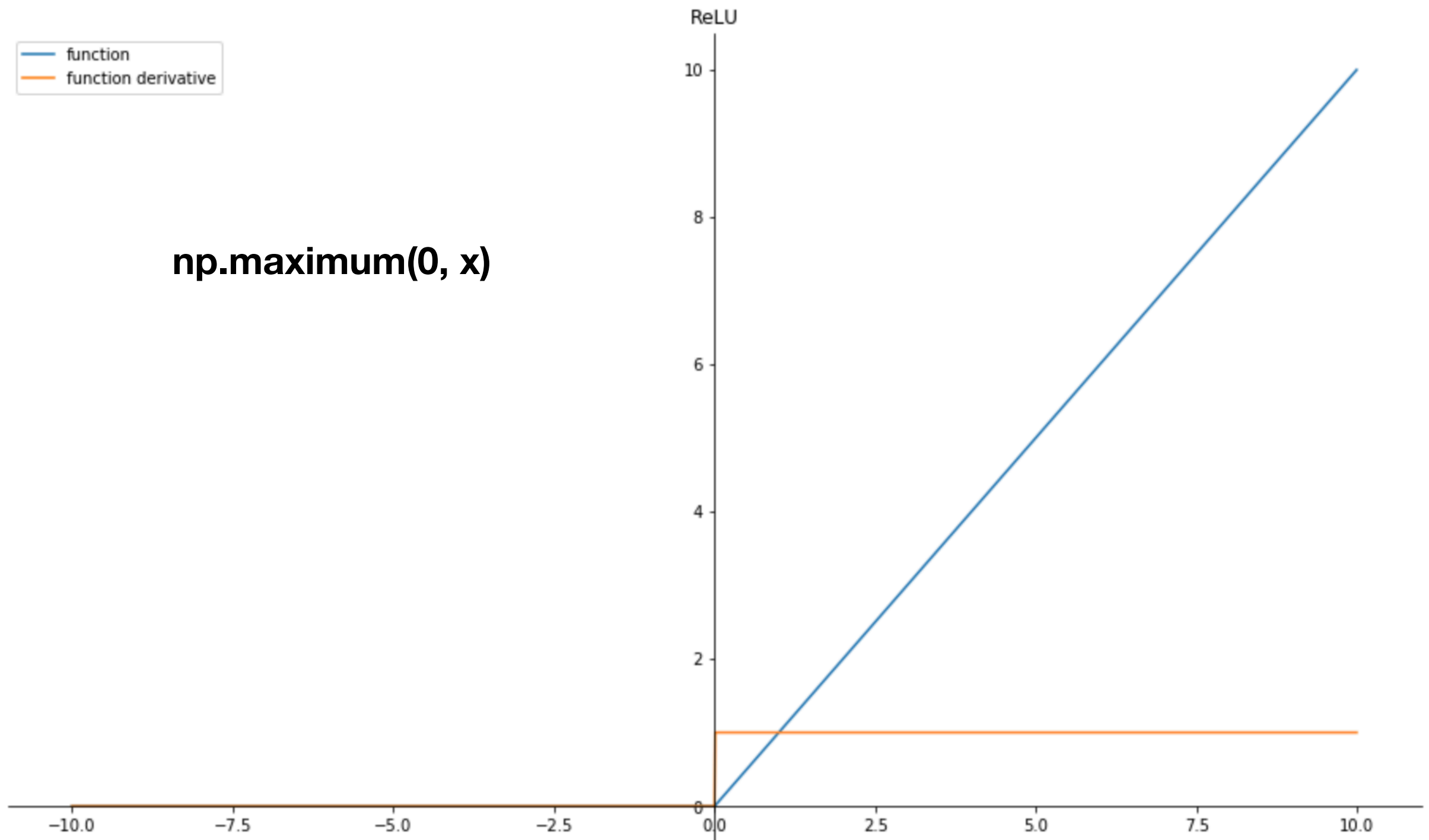
Sigmoid



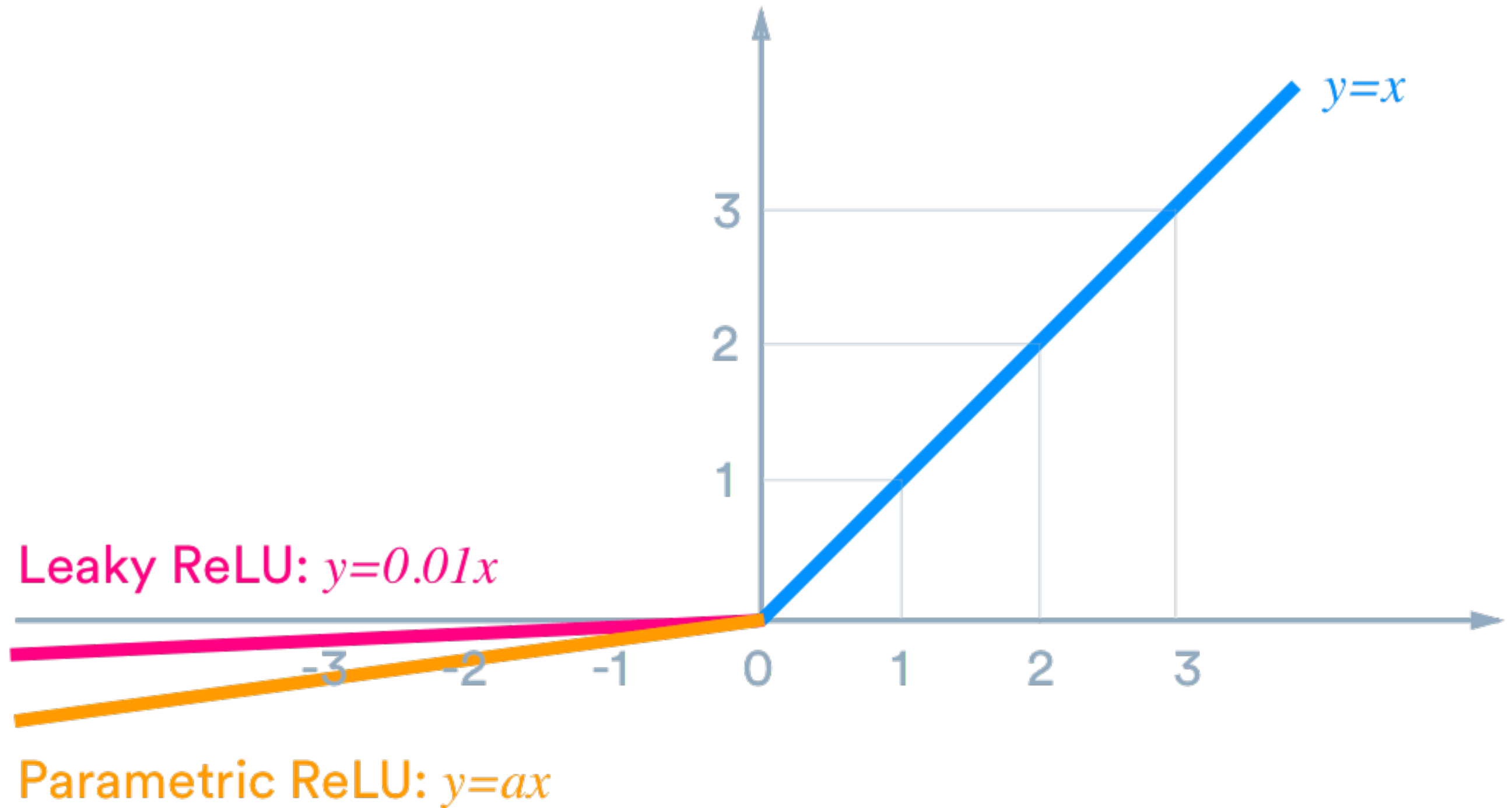
Tanh



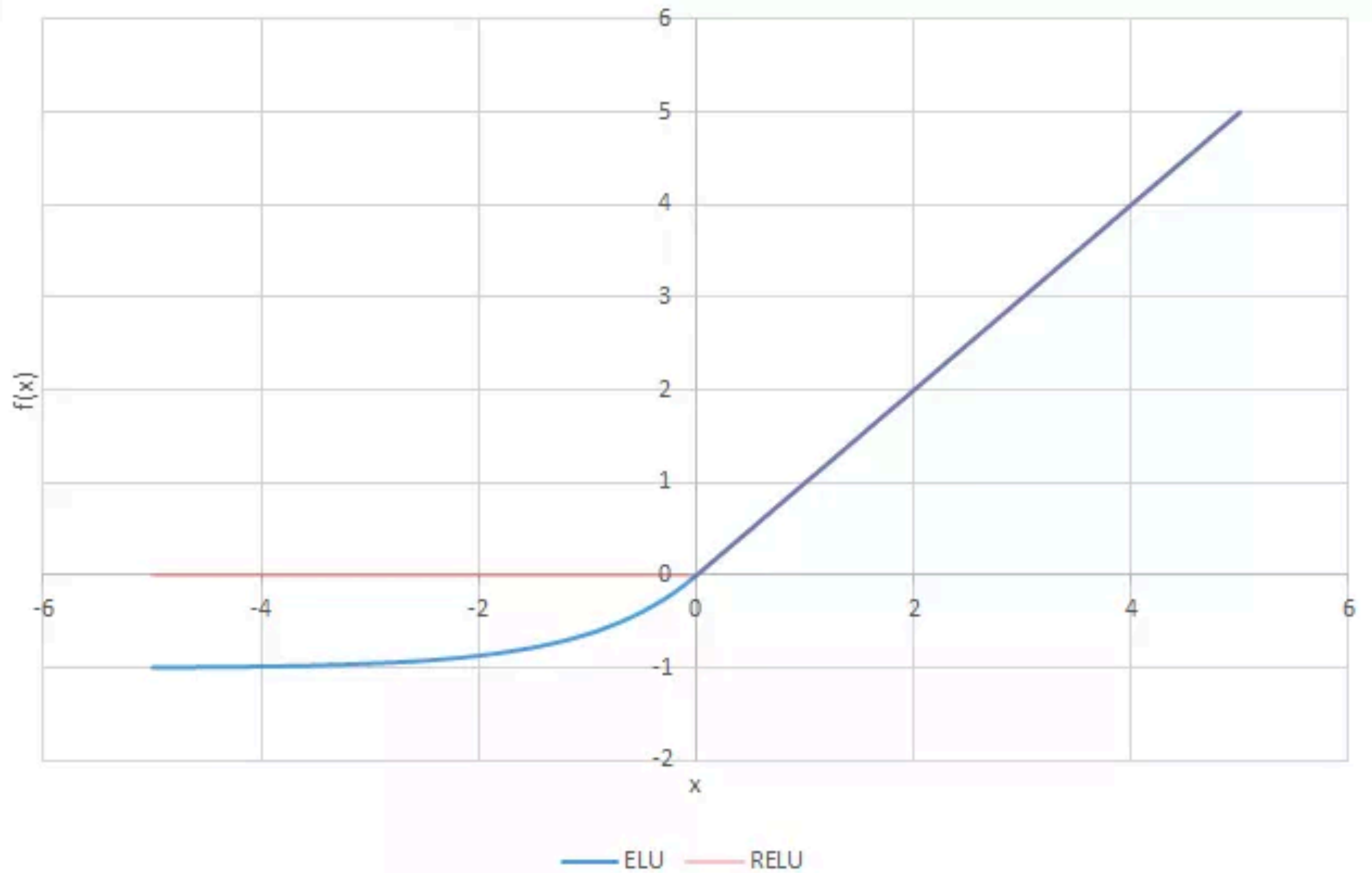
ReLU



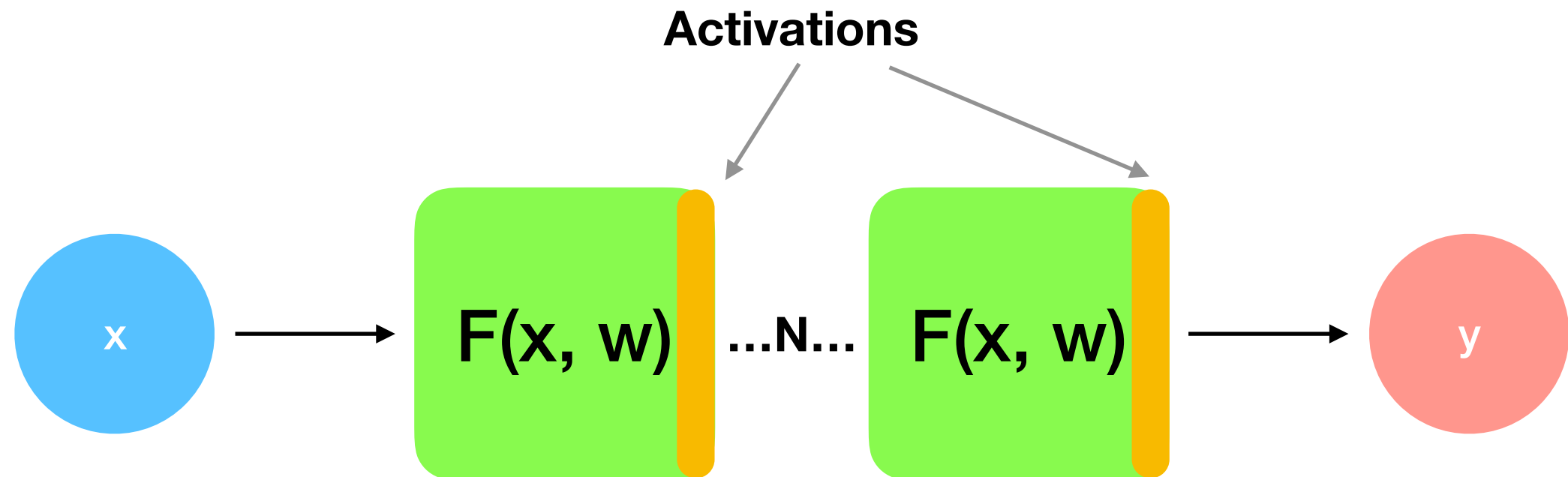
Leaky ReLU



ELU



Neural Networks

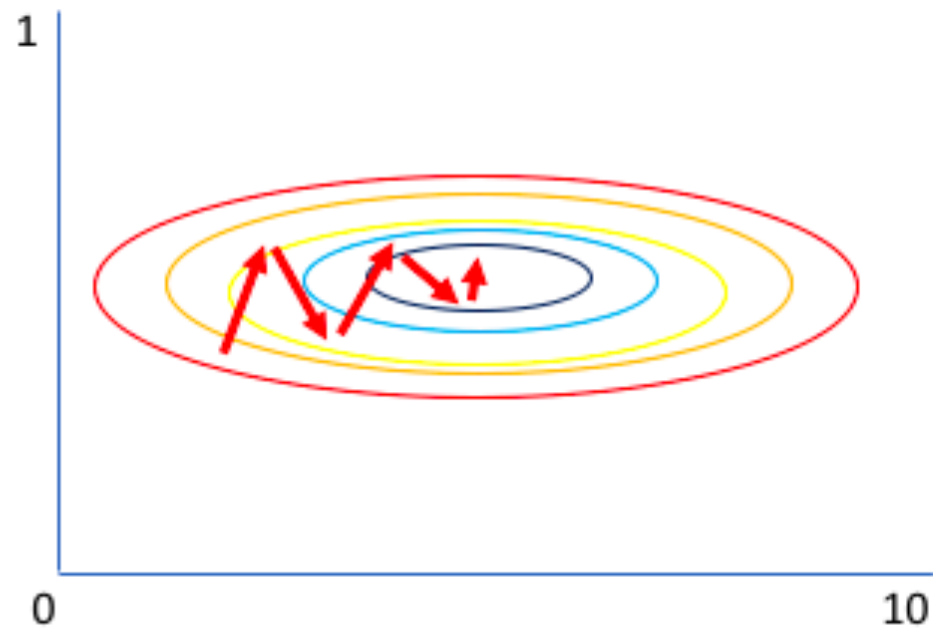


Important things

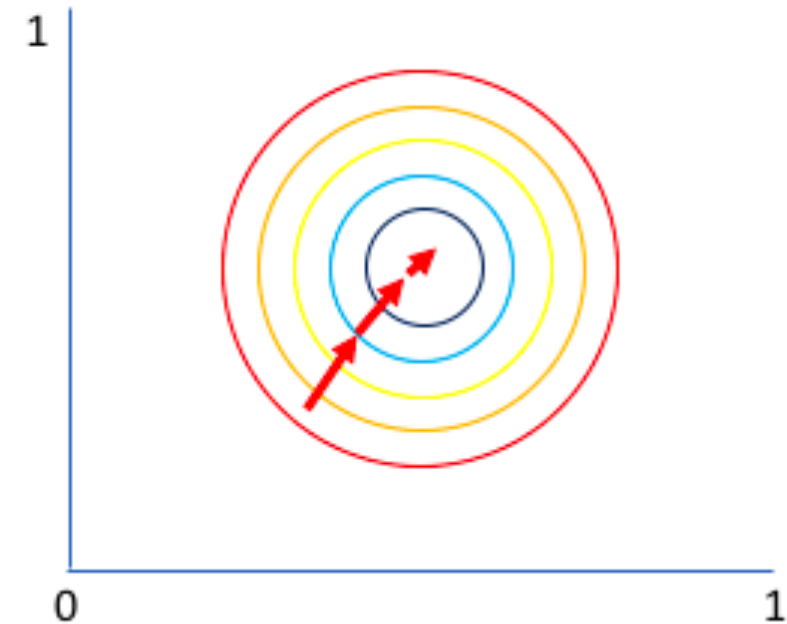
- Don't use sigmoid as inner activation in MLP
- But then we talk about how you can use sigmoid as inner activation
- Normalize your data
- Early stopping

Important things

Why normalize?



Gradient of larger parameter
dominates the update



Both parameters can be
updated in equal proportions

Early Stopping

