

DEEP NEURAL NETWORK

Deep Learning for Computer Vision

Arthur Douillard

Deep Neural Networks

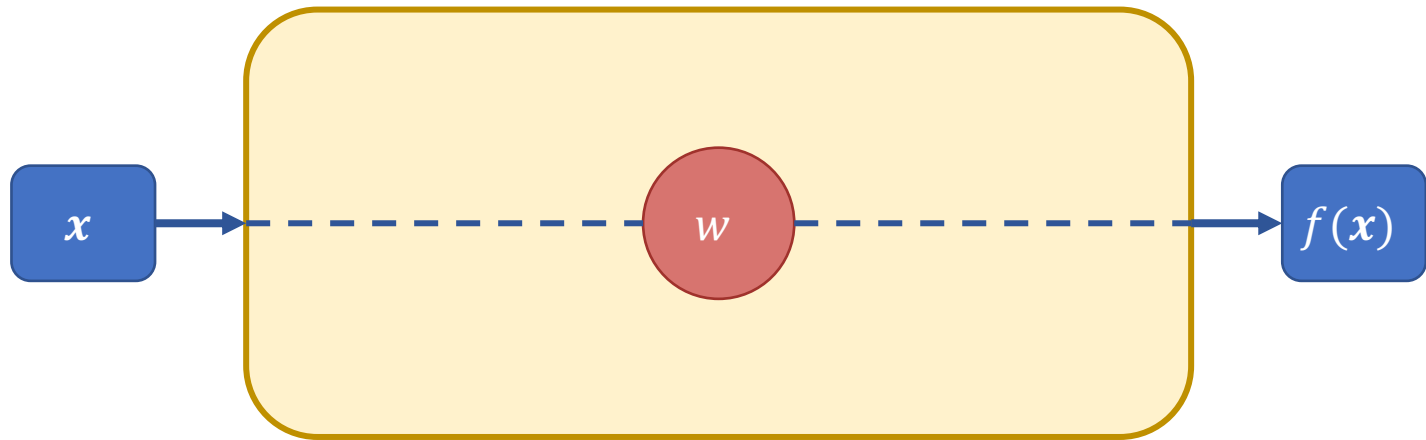


A neural network approximates a function





Linear regression

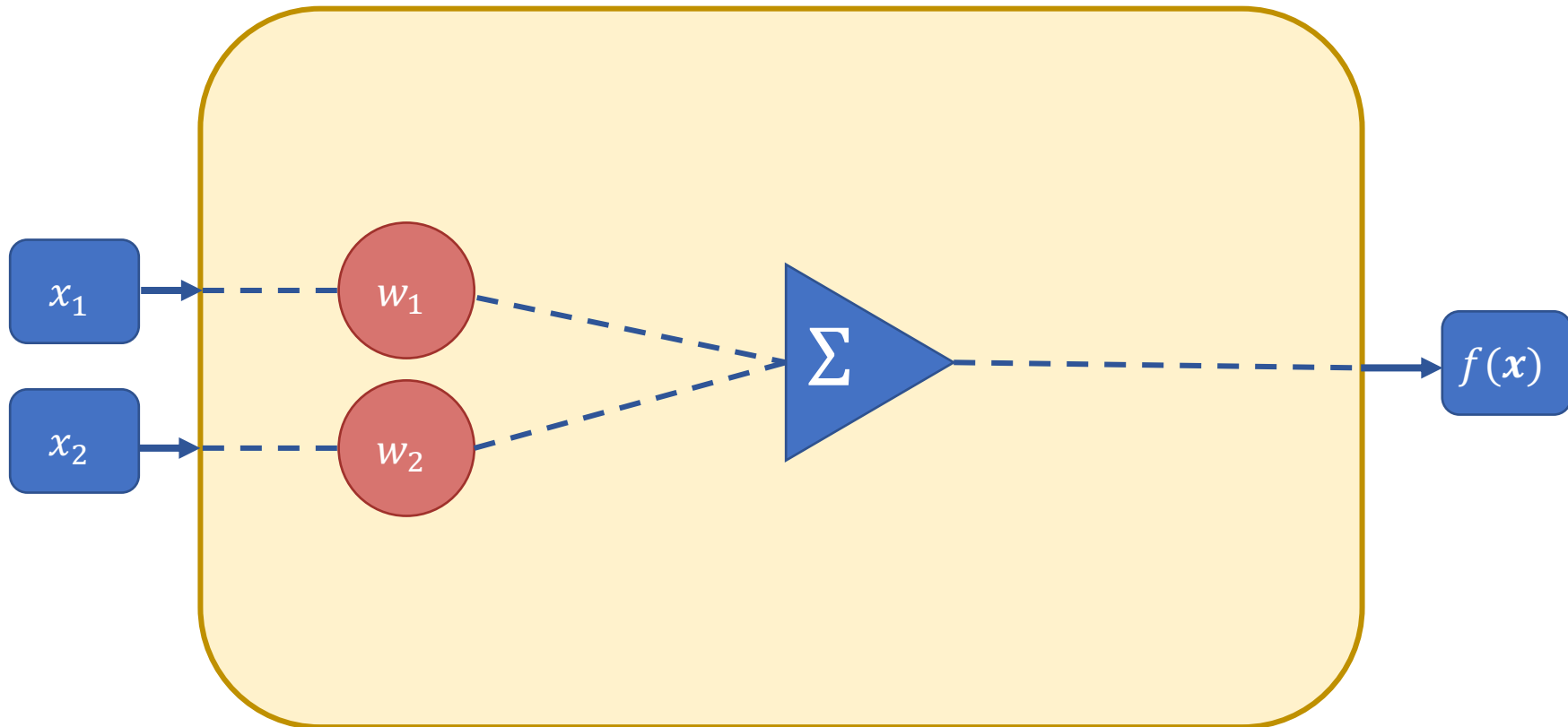


Useful to predict continuous variables

$$f(\mathbf{x}) = \mathbf{w}\mathbf{x}$$



Linear regression

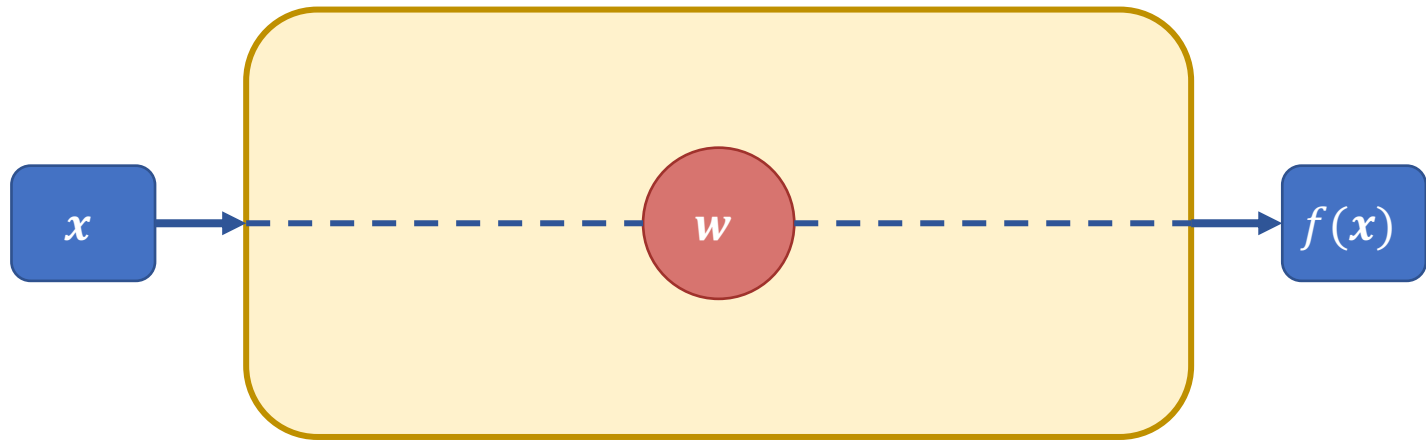


The “connection” often depicted are only multiplications and additions

$$f(\mathbf{x}) = w_1 x_1 + w_2 x_2$$



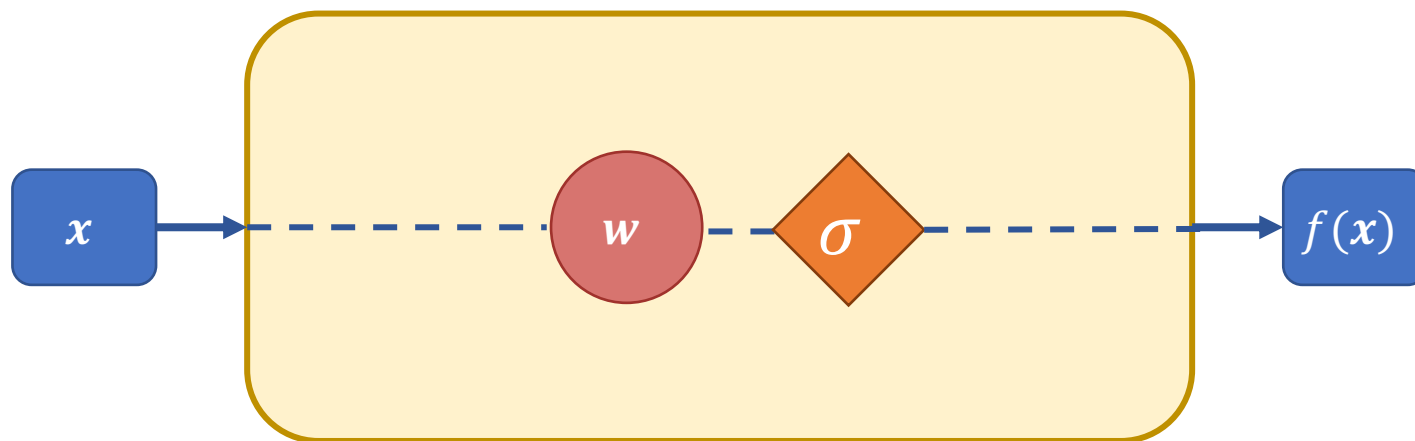
Linear regression



$$f(x) = wx$$



Single Neuron



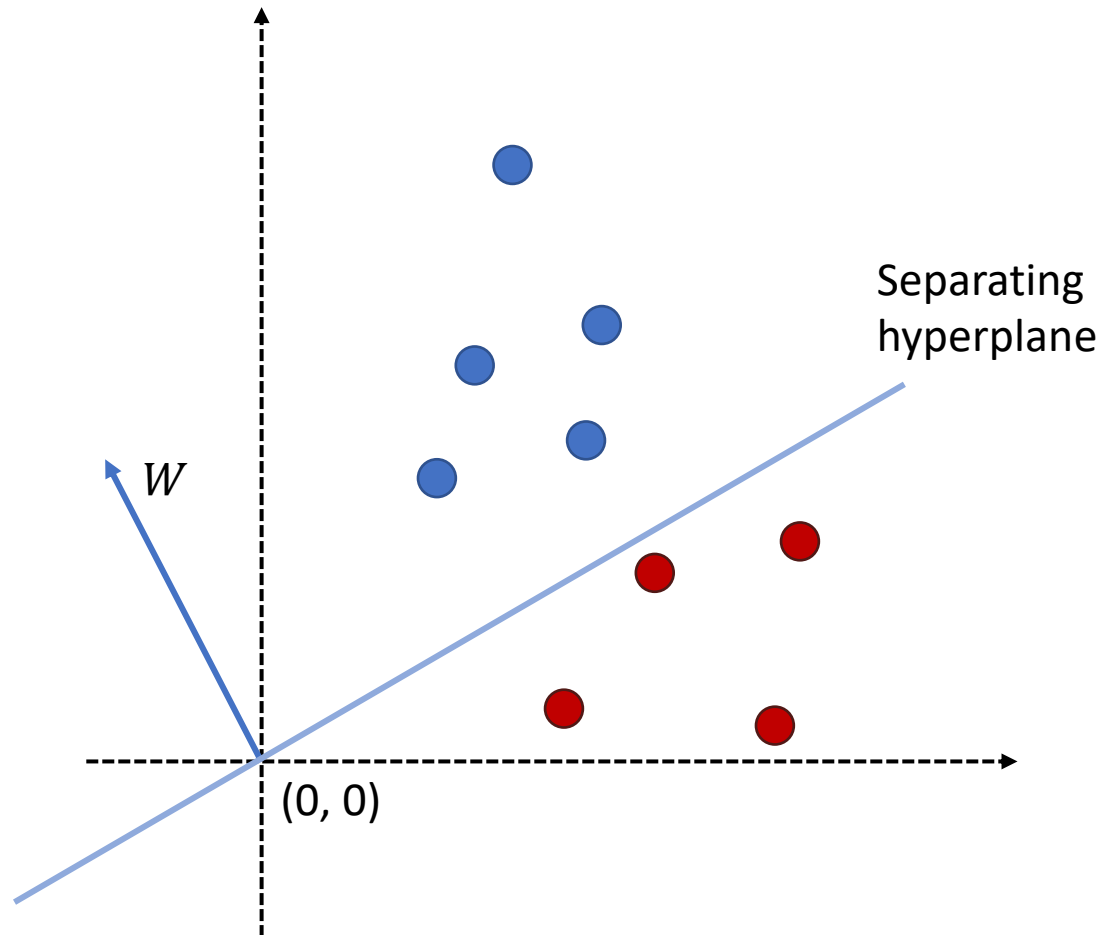
Non-linear activation to determine how much a neuron “fires”

$$f(\mathbf{x}) = \sigma(\mathbf{w}\mathbf{x})$$

Separating Hyperplane

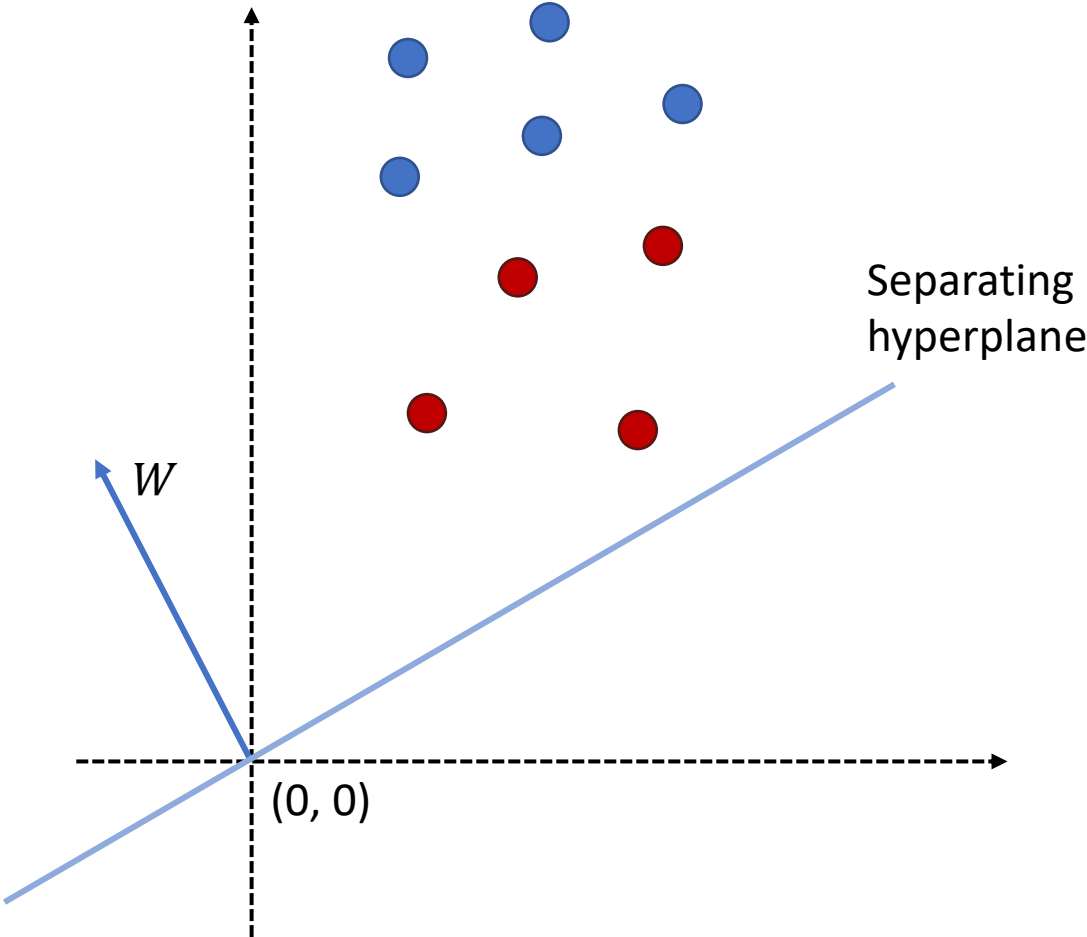


Optimize so that W is orthogonal to the separating hyperplane



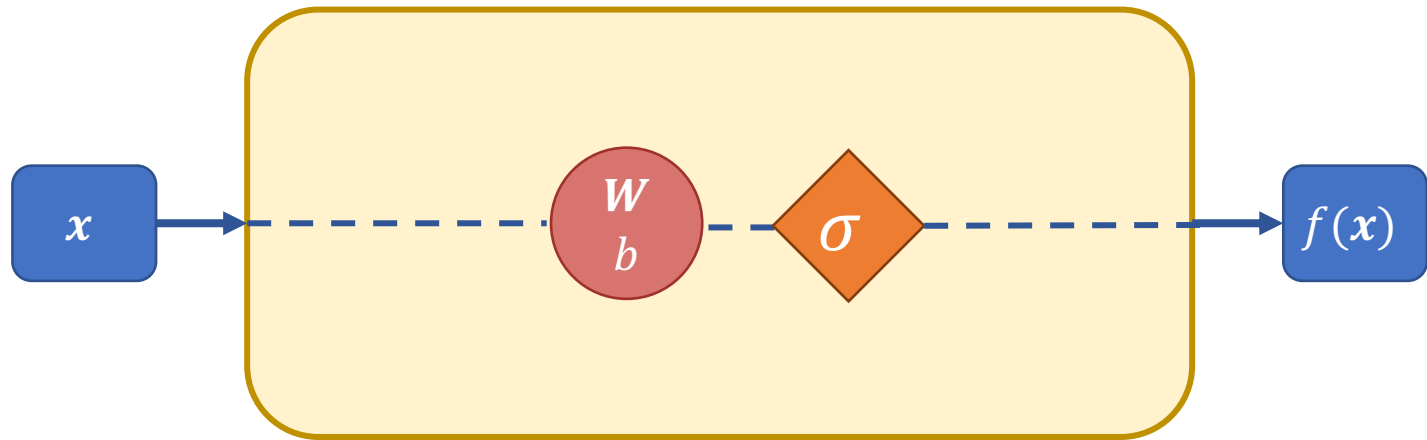


Need bias!



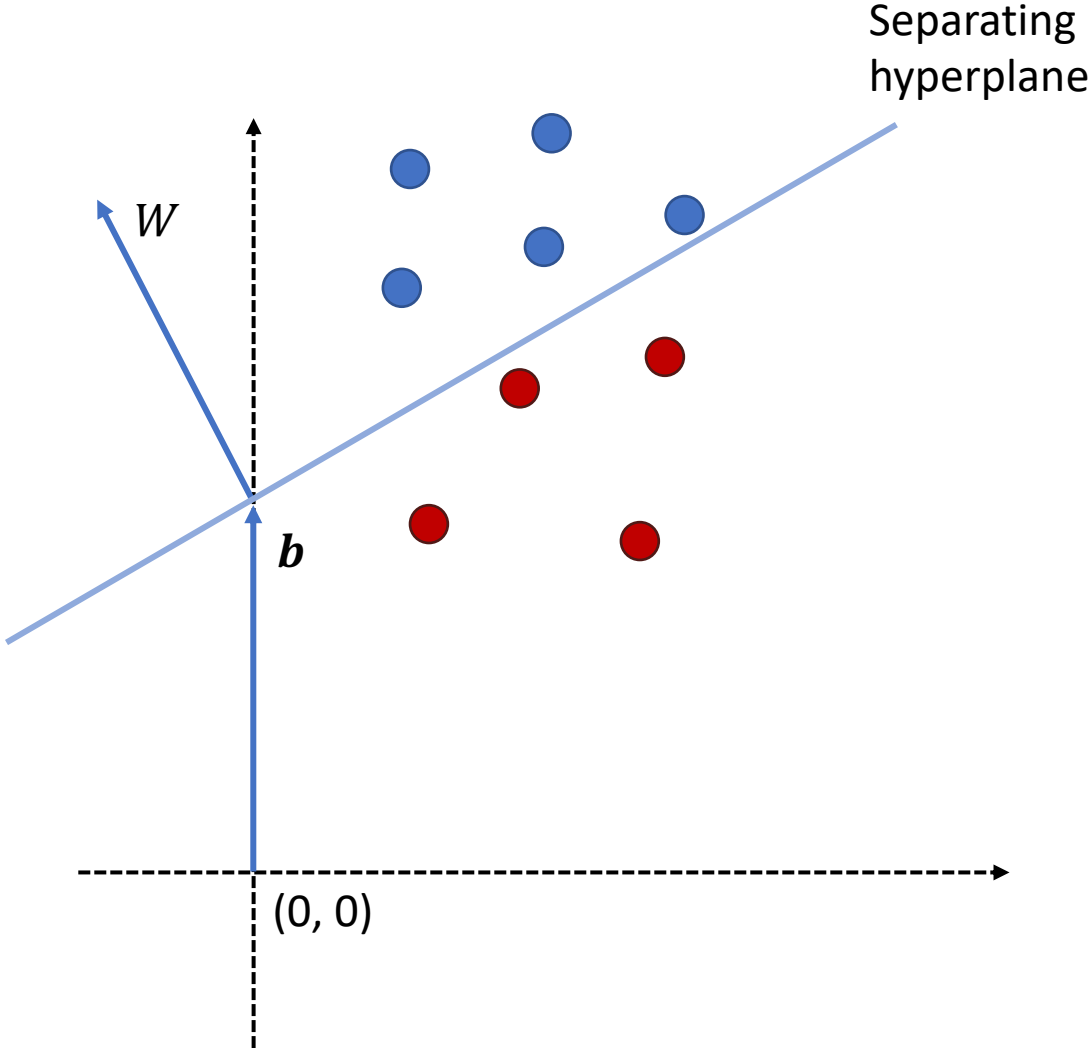


Single Neuron



With a **bias**

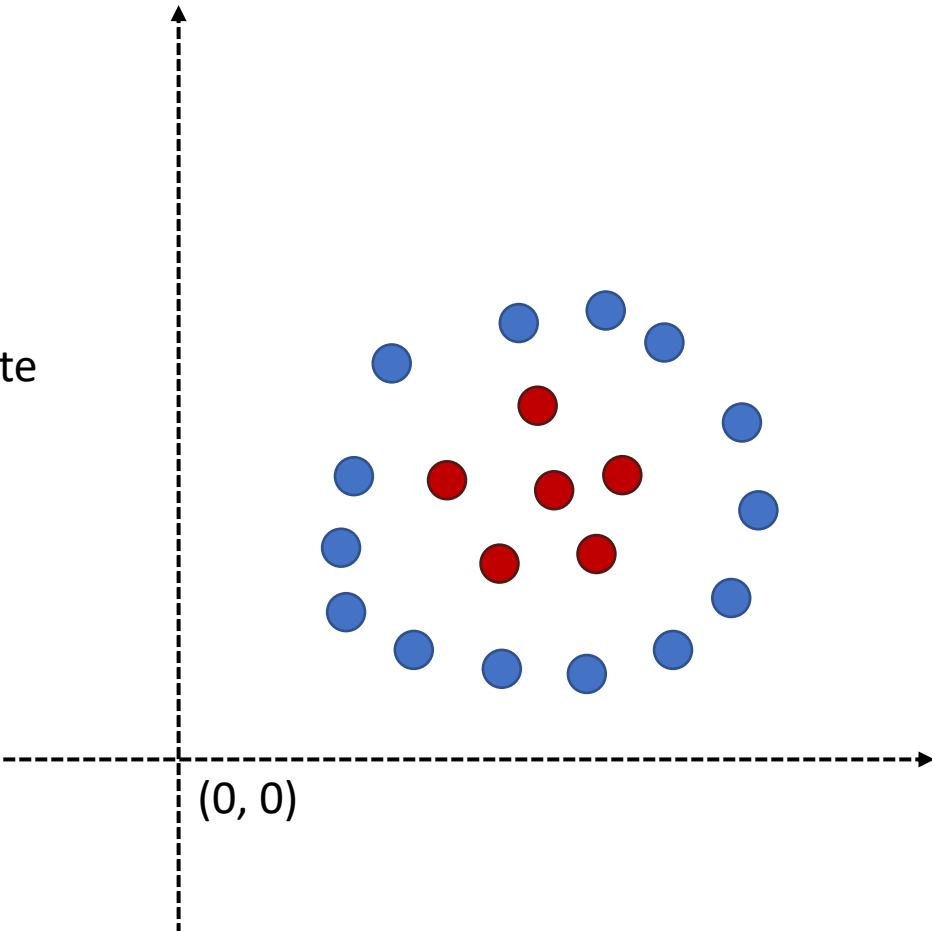
$$f(x) = \sigma(wx + b)$$



Non-Linear Data

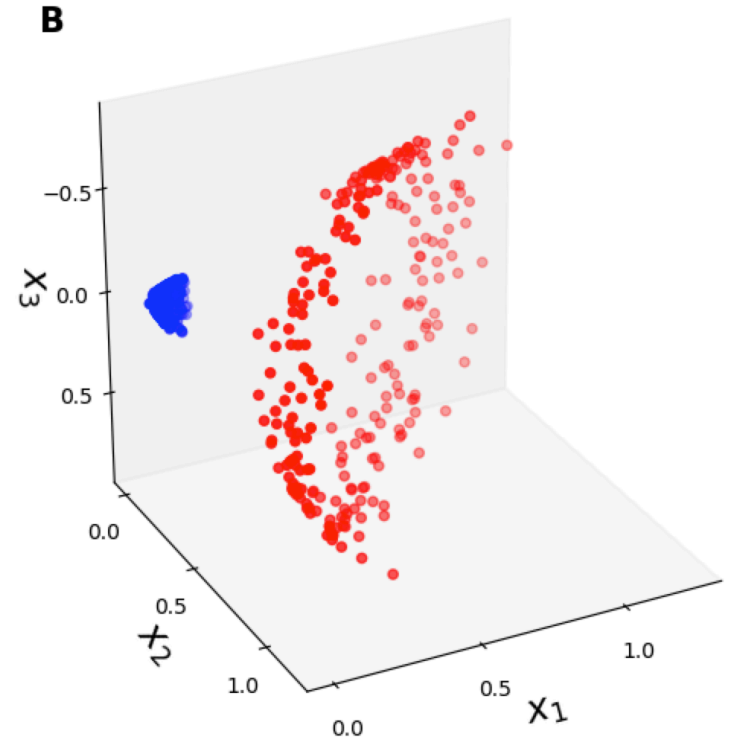
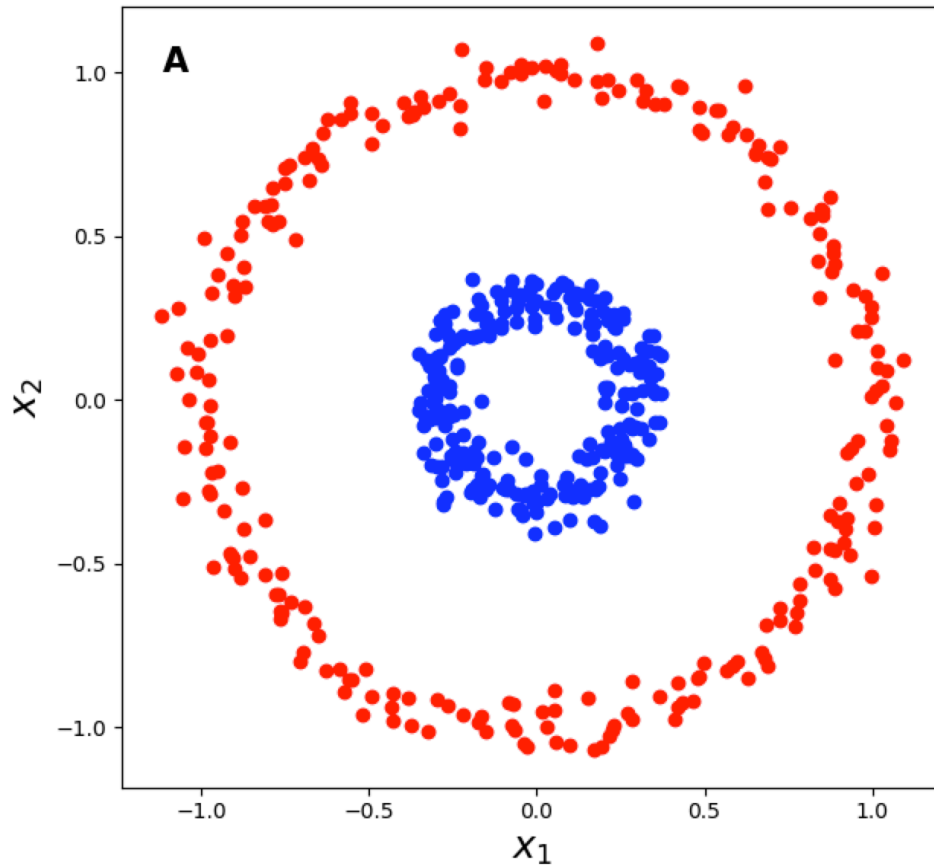


A linear model cannot discriminate **blue** and **red** classes!





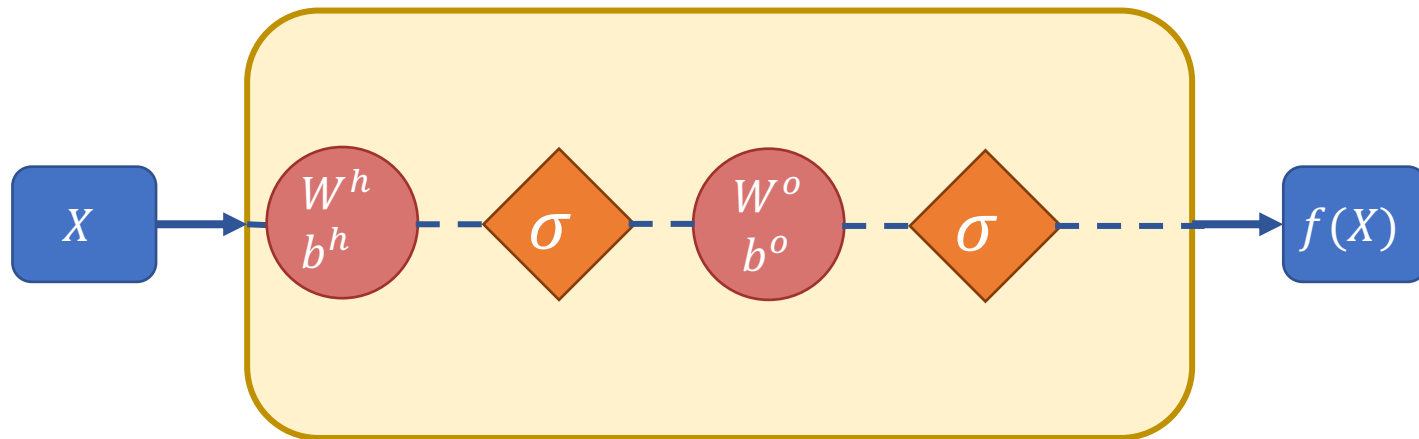
Before, you would have use the **kernel trick**



Excellent blog post on that [Gundersen](#)



Multi-Layer Perceptron



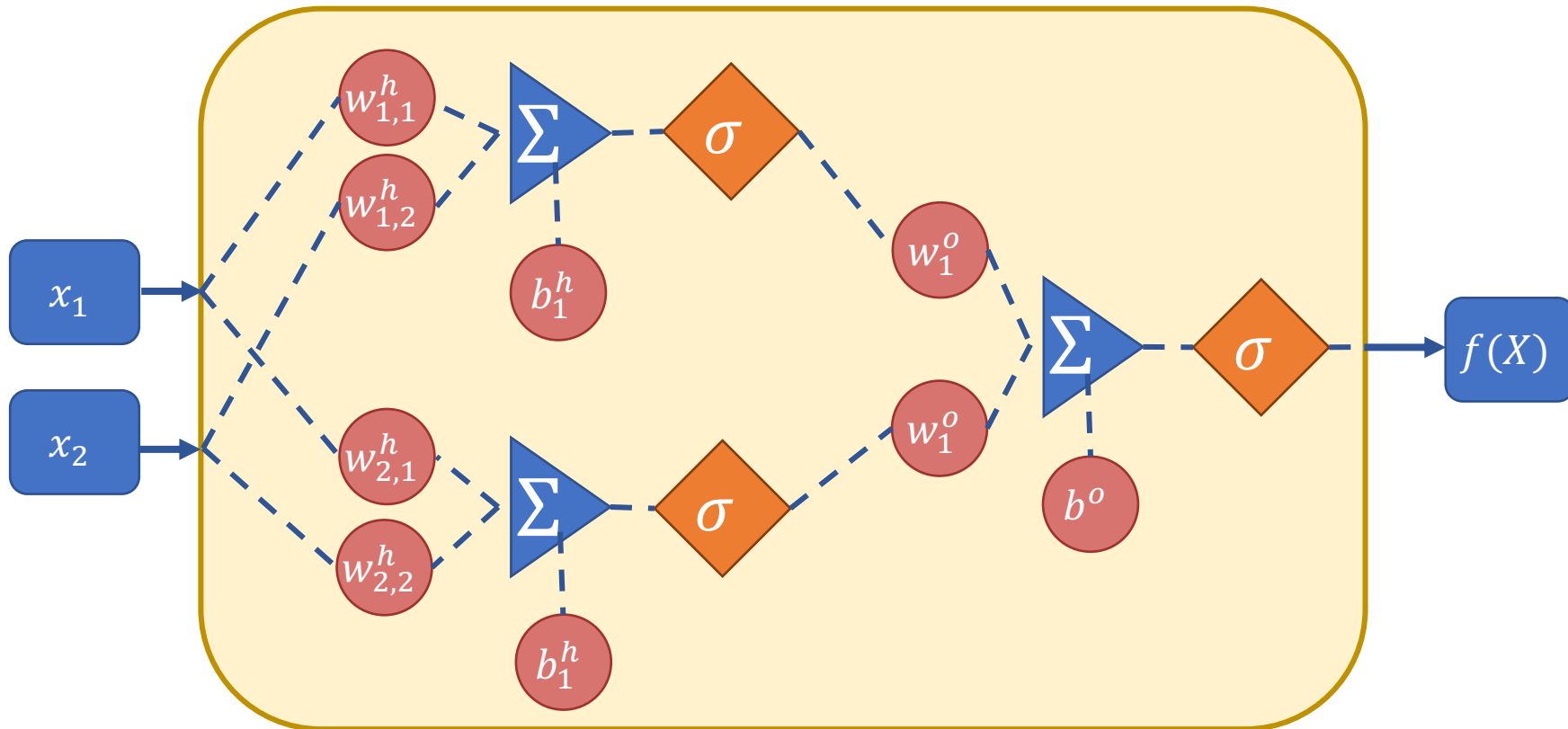
Stack multiple **linear transformations** and **non-linear activation**

$$f(\mathbf{x}) = \sigma(\mathbf{w}^o \sigma(\mathbf{W}^h \mathbf{x} + \mathbf{b}^h) + \mathbf{b}^o)$$

Multi-Layer Perceptron

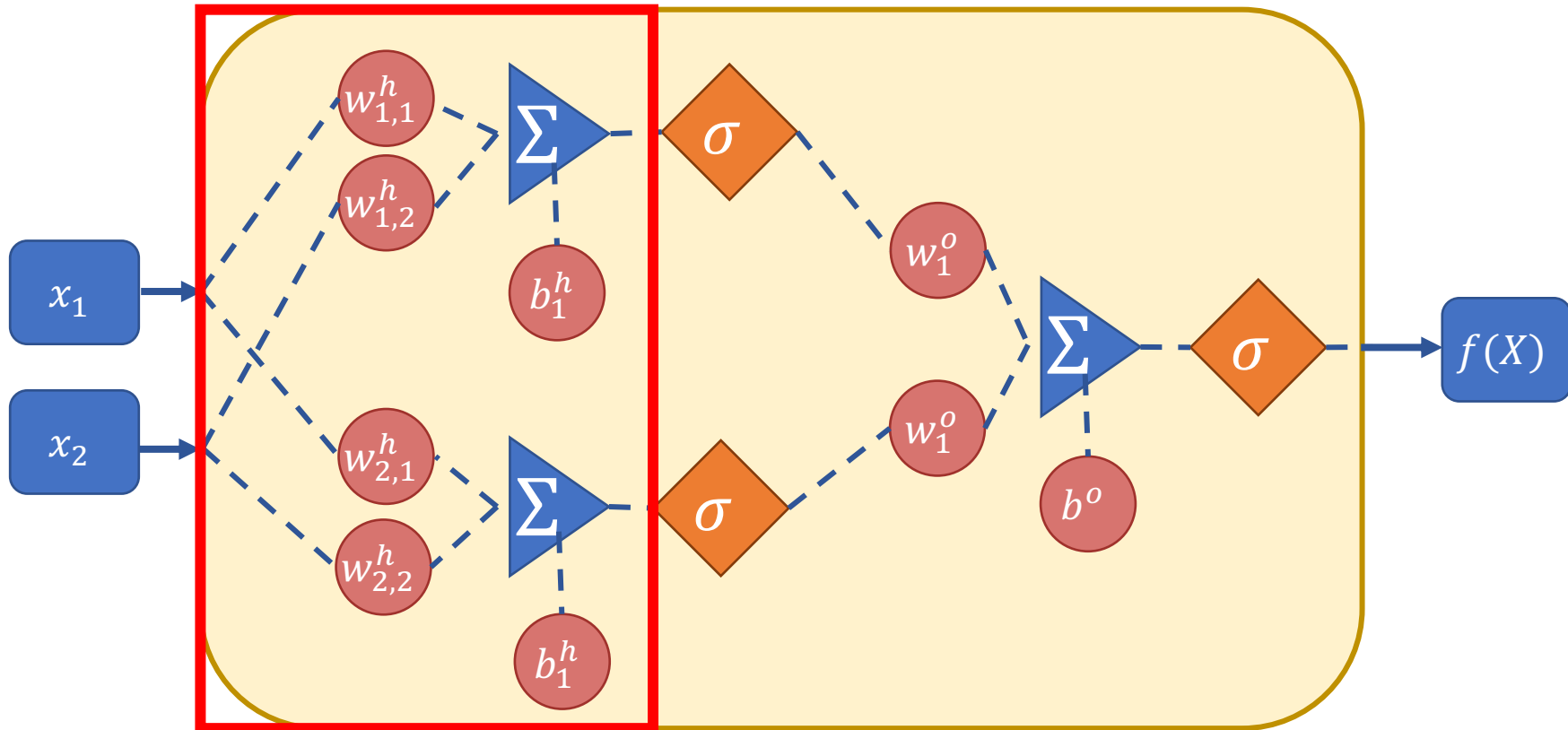


W^h is a matrix because it has here 2 output dimensions



$$f(\mathbf{x}) = \sigma(\mathbf{w}^o \sigma(\mathbf{W}^h \mathbf{x} + \mathbf{b}^h) + b^o)$$

Multi-Layer Perceptron

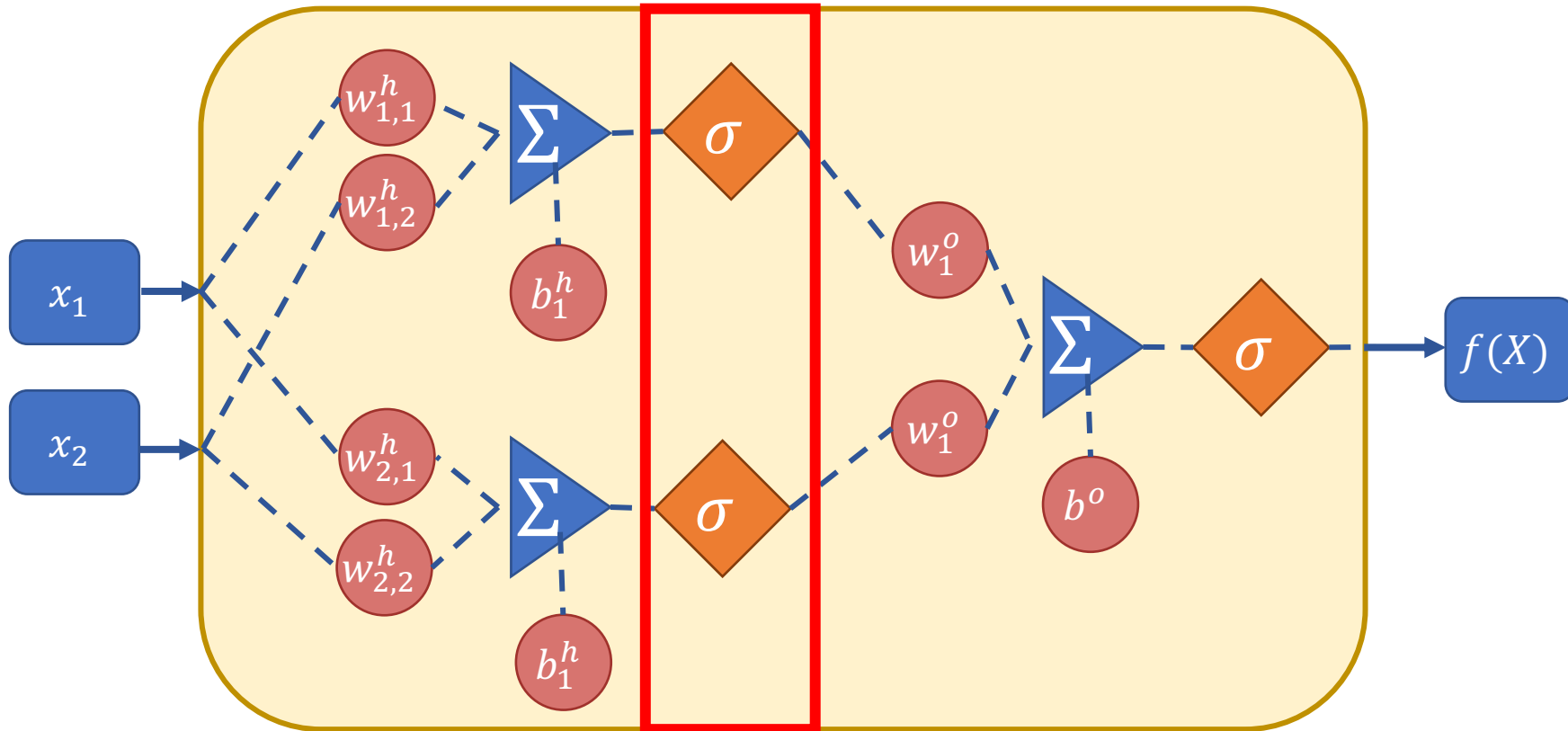


$$\tilde{\mathbf{h}} = \mathbf{W}^h \mathbf{x} + \mathbf{b}^h$$

Multi-Layer Perceptron



Embeddings/features

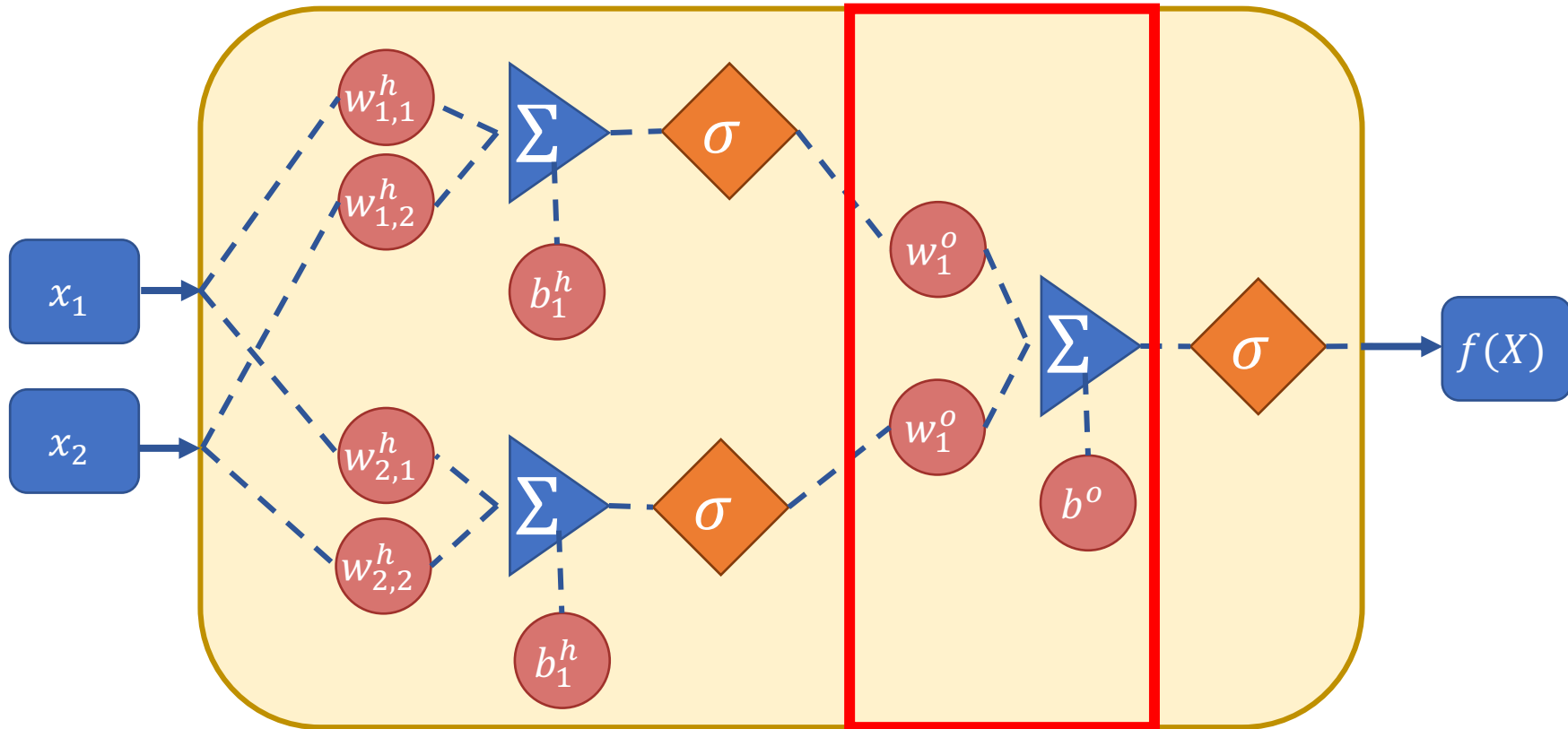


$$\mathbf{h} = \sigma(\tilde{\mathbf{h}})$$

Multi-Layer Perceptron

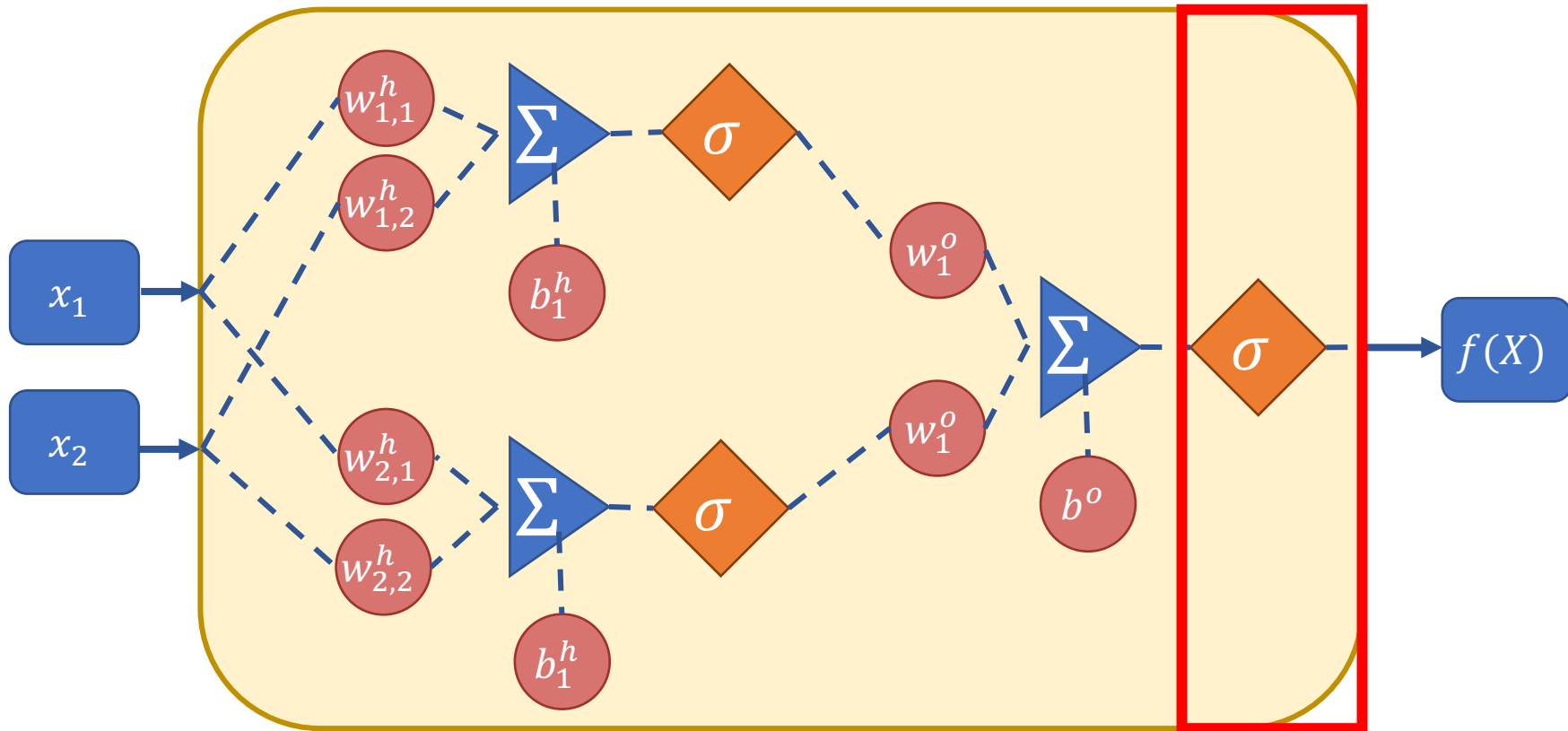


Logits!



$$\tilde{y} = w^o h + b^o$$

Multi-Layer Perceptron

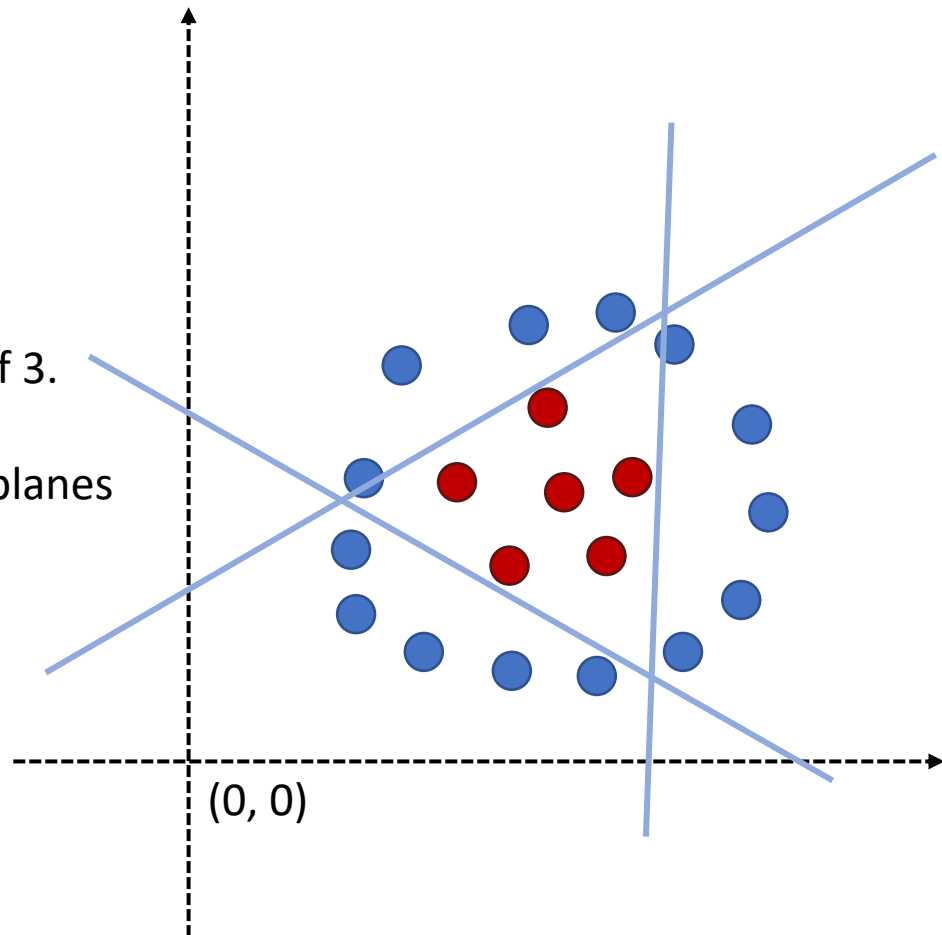


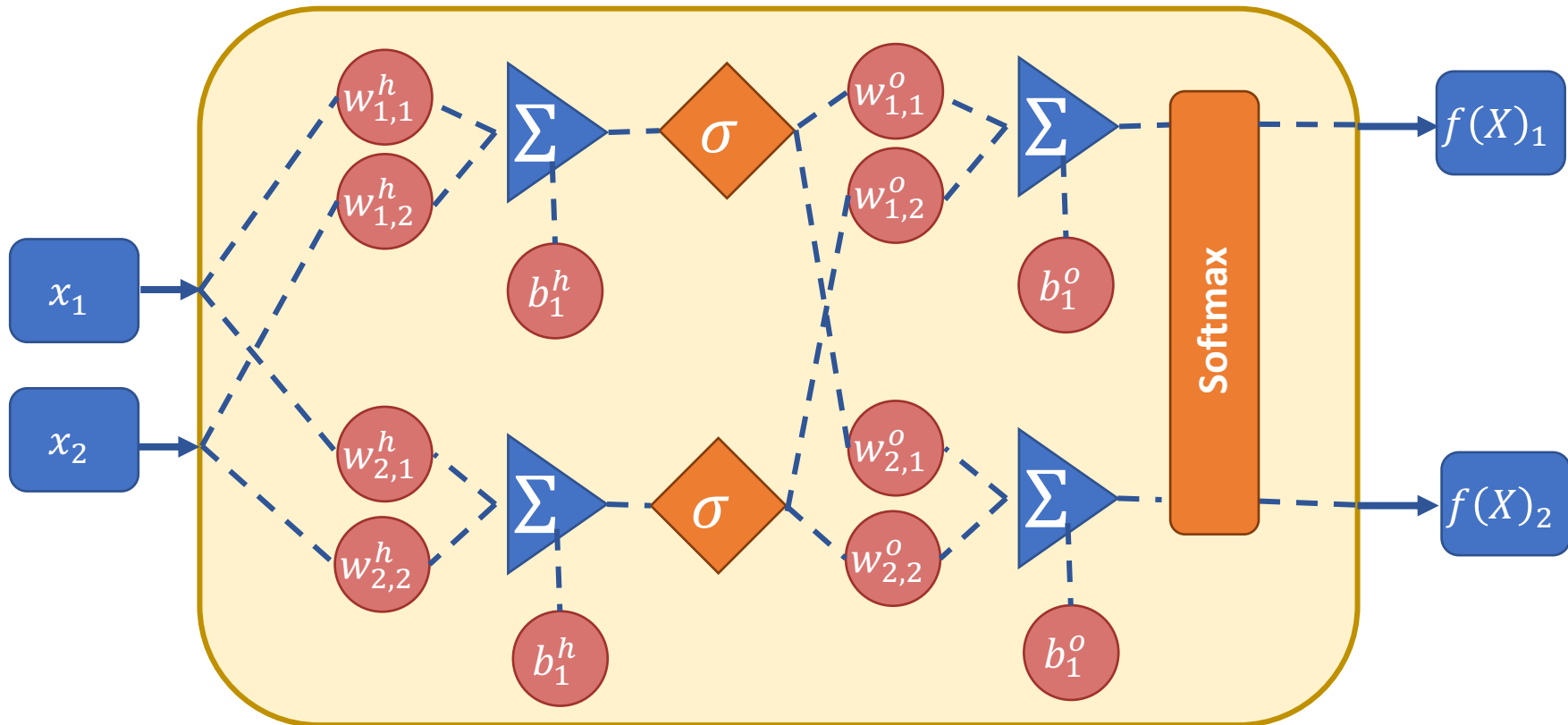
$$\hat{\mathbf{y}} = \sigma(\tilde{\mathbf{y}})$$



Considering hidden dimension of 3.

Hidden layer can define 3 hyperplanes

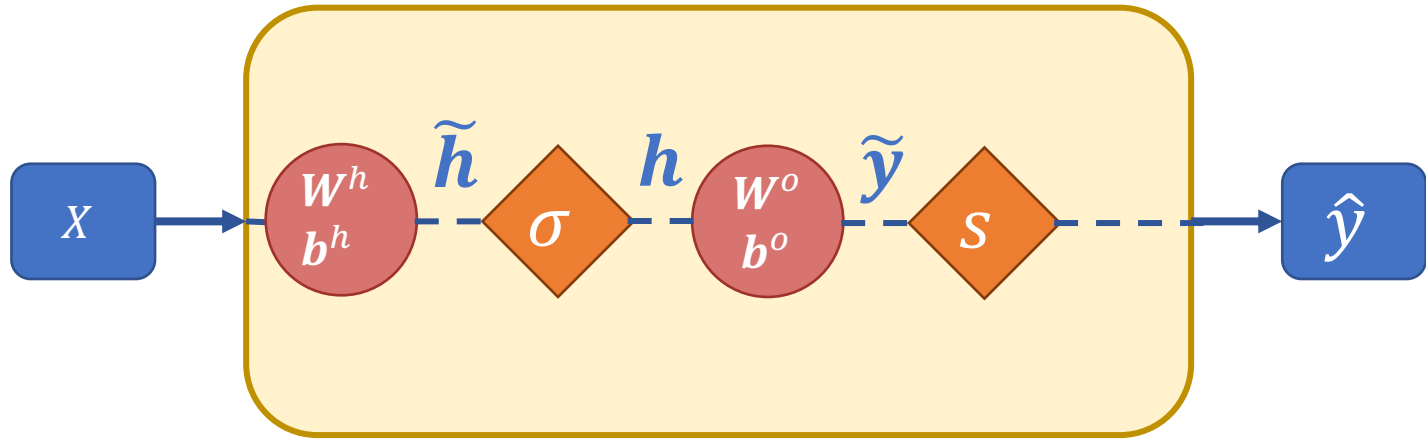




W^o is a matrix because it has here 2 predicted classes.
Can be extended to 3, 4, ..., 1000 classes.



Multi-Layer Perceptron



$$\tilde{\mathbf{h}} = \mathbf{W}^h \mathbf{x} + \mathbf{b}^h$$

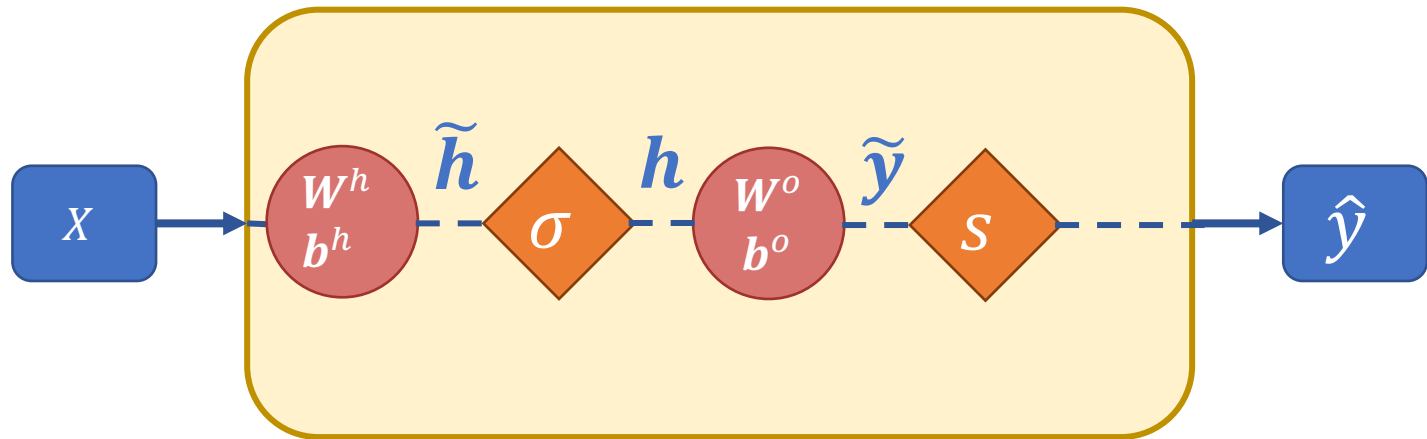
$$\mathbf{h} = \sigma(\tilde{\mathbf{h}})$$

$$\tilde{\mathbf{y}} = \mathbf{W}^o \mathbf{h} + \mathbf{b}^o$$

$$\hat{\mathbf{y}} = \sigma(\tilde{\mathbf{y}})$$



Multi-Layer Perceptron



$$\tilde{h} = W^h x + b^h$$

$$h = \sigma(\tilde{h})$$

$$\tilde{y} = W^o h + b^o$$

$$\hat{y} = \sigma(\tilde{y})$$

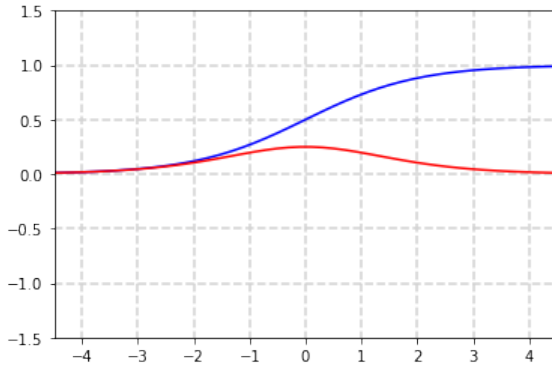
→ features / embeddings

→ logits

→ Model predictions



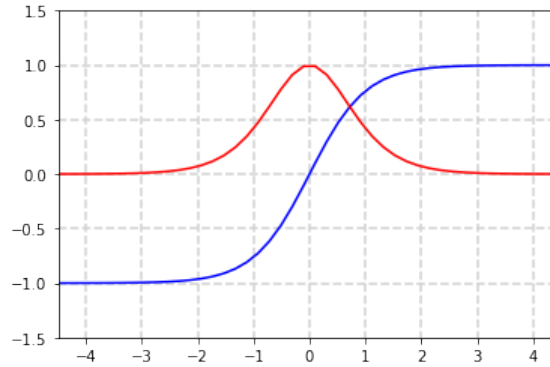
Function and their derivative



Sigmoid

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

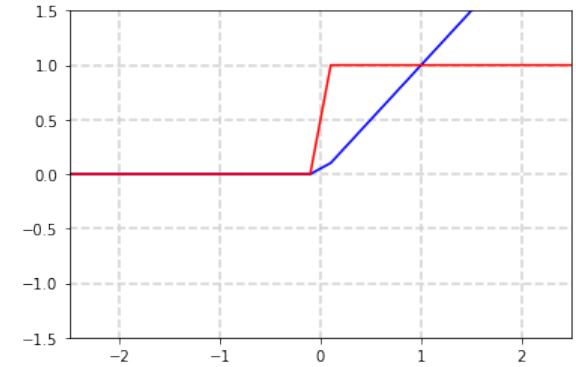
$$\frac{d}{dx} \sigma(x) = \sigma(x)(1 - \sigma(x))$$



Hyperbolic tangent

$$\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$$

$$\frac{d}{dx} \tanh(x) = 1 - \tanh(x)^2$$



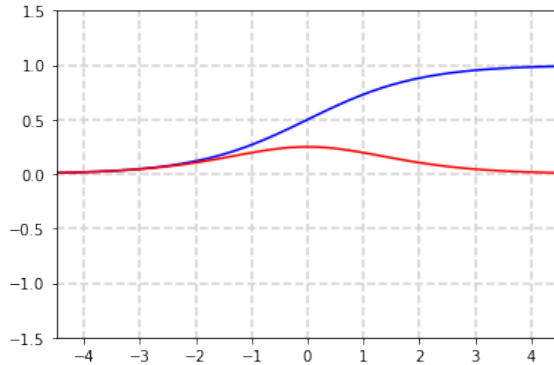
Rectified Linear Unit

$$\text{ReLU}(x) = \max(x, 0)$$

$$\frac{d}{dx} \text{ReLU}(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$



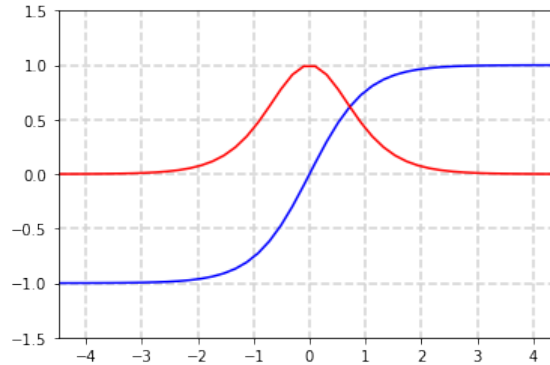
Function and their derivative



Sigmoid

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

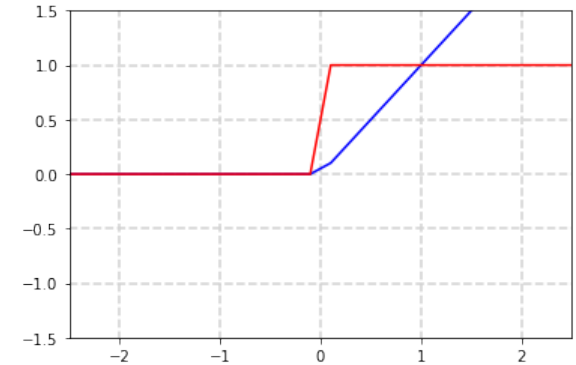
$$\frac{d}{dx} \sigma(x) = \sigma(x)(1 - \sigma(x))$$



Hyperbolic tangent

$$\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$$

$$\frac{d}{dx} \tanh(x) = 1 - \tanh(x)^2$$



Rectified Linear Unit

$$\text{ReLU}(x) = \max(x, 0)$$

$$\frac{d}{dx} \text{ReLU}(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

Without hidden activation, a MLP is equivalent to a single layer!

→ Composition of affine functions is an affine function



Binary classification: **sigmoid**

→ Applied element-wise

→ Range [0, 1]

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\frac{d}{dx}\sigma(x) = \sigma(x)(1 - \sigma(x))$$

Multi-Class classification: **softmax**

→ Applied over a vector

→ Range [0, 1] per element

→ Sum to 1 for the vector

→ Probability distribution

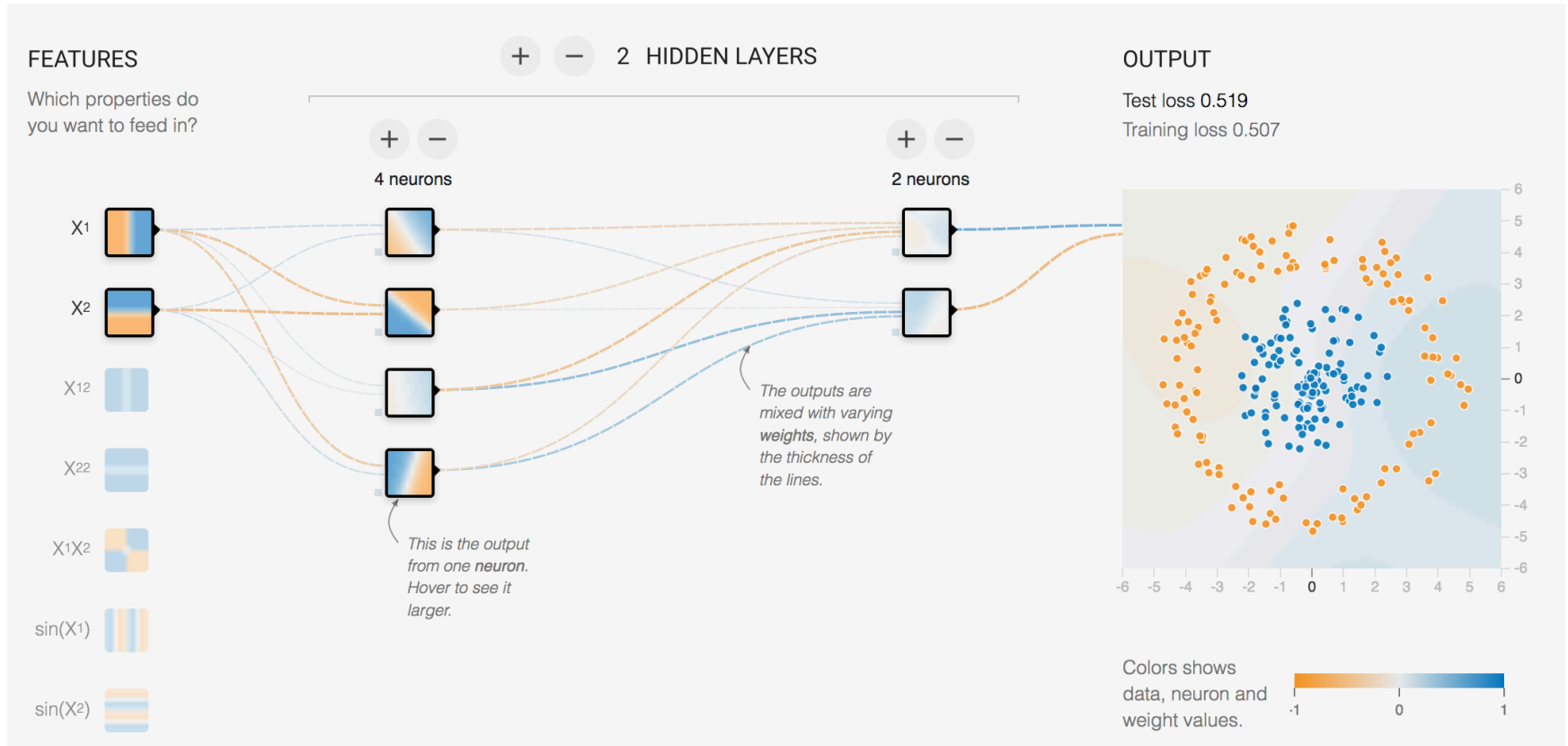
$$\text{softmax}(\mathbf{x}) = \frac{1}{\sum_j e^{x_j}} \begin{bmatrix} e^{x_1} \\ \dots \\ e^{x_n} \end{bmatrix}$$

$$\text{softmax}(\mathbf{x})_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

$$\frac{d}{dx_j}\text{softmax}(\mathbf{x})_i = \begin{cases} \text{softmax}(\mathbf{x})_i(1 - \text{softmax}(\mathbf{x})_i) & \text{if } i = j \\ -\text{softmax}(\mathbf{x})_i\text{softmax}(\mathbf{x})_j & \text{if } i \neq j \end{cases}$$



To get an intuition of the effect of hidden dimensions, number of layers, and activations:

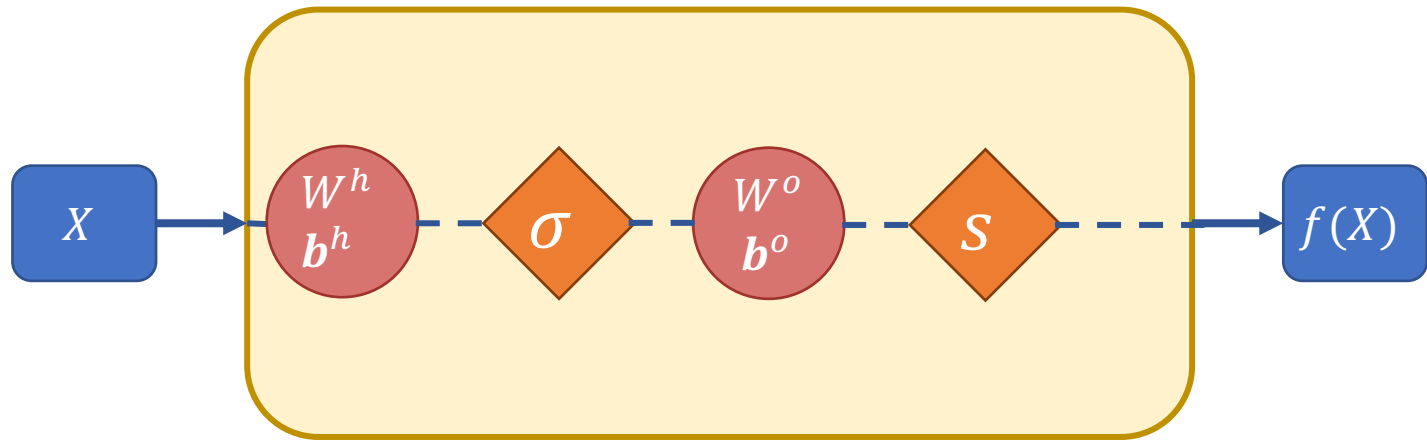


playground.tensorflow.org

Learning DNNs

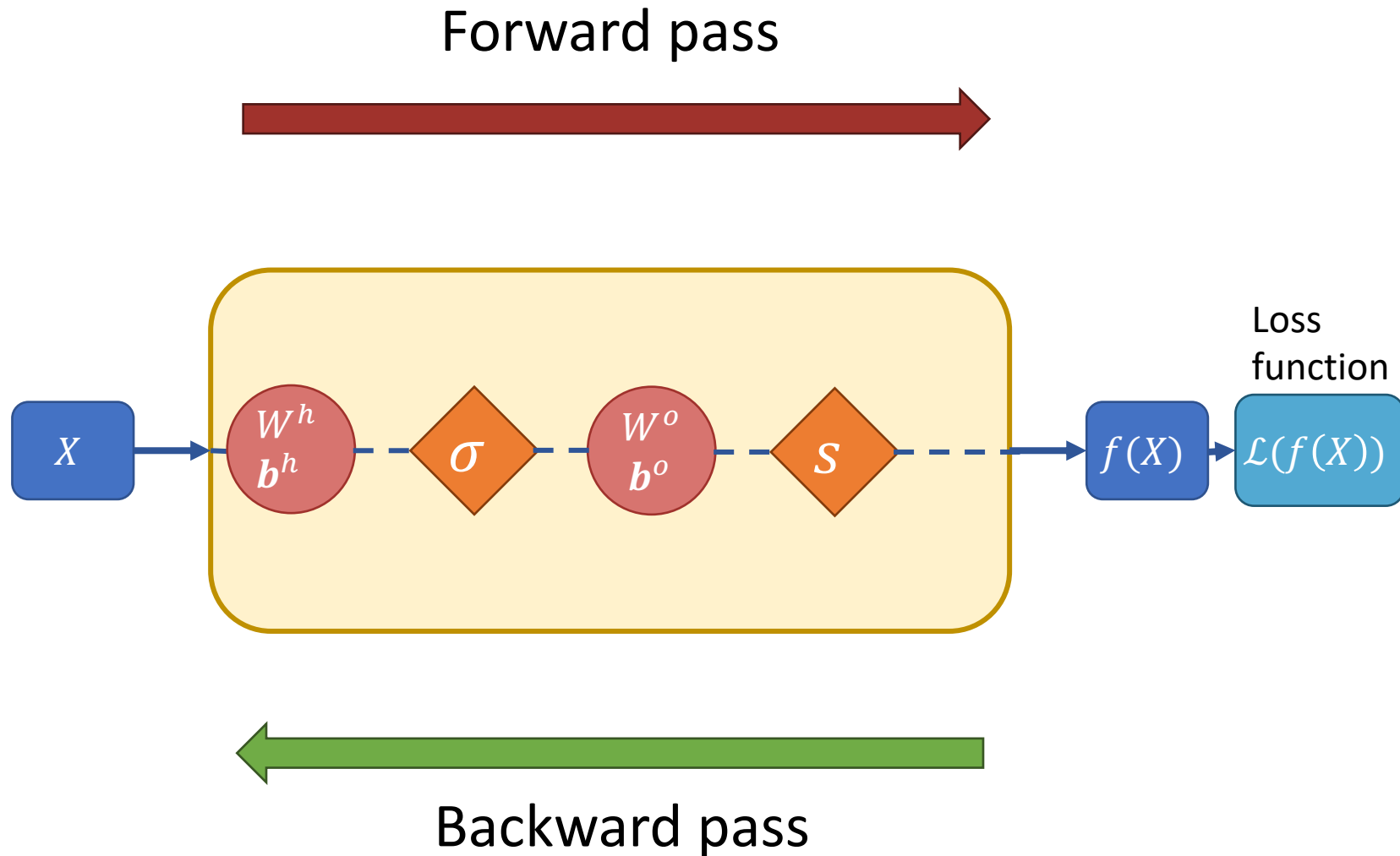


Multi-Layer Perceptron



1 hidden layer, H hidden dimensions, C output dimensions with a softmax

Forward & Backward passes





For classification with softmax as final activation:

Cross-entropy (also known as negative log-likelihood):

$$\mathcal{L}_{CE}(\hat{y}, y) = - \sum_i y_i \log \hat{y}_i$$

One-hot
target



Dog: 0.2

Cat: 0.8

Cross-Entropy: $-\log 0.8 = 0.2231 \dots$



Avoid doing *if* in GPUs, use one-hot in cross-entropy.

Given 5 classes, if the ground-truth class is 3:

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Zero-indexed!



Optimize the network to minimize the loss with respect to all neurons:

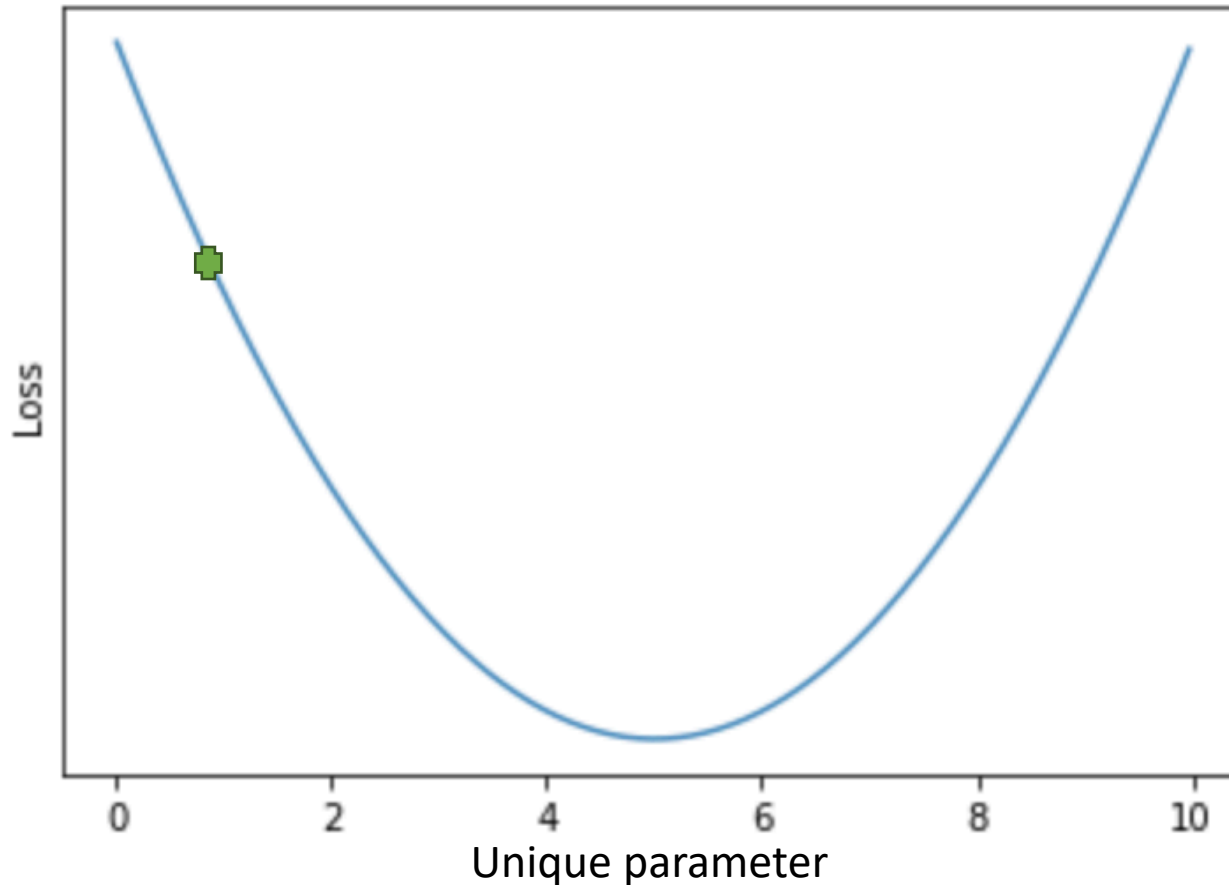
$$\mathcal{L}_{CE}(\hat{y}, y) = - \sum_i y_i \log \hat{y}_i$$

Minimization of a function



Optimize the network to minimize the loss with respect to all neurons:

$$\mathcal{L}_{CE}(\hat{y}, y) = - \sum y_i \log \hat{y}_i$$

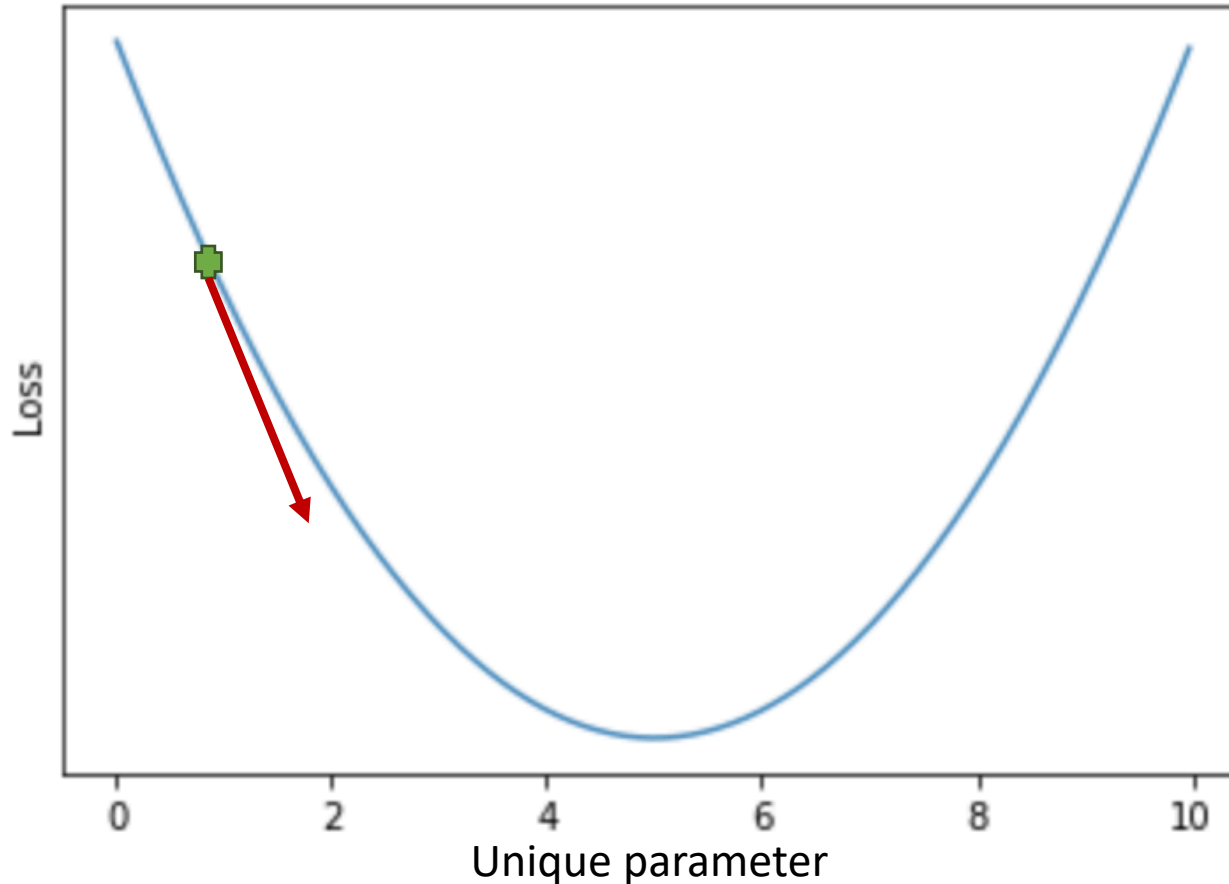


Minimization of a function



Optimize the network to minimize the loss with respect to all neurons θ :

$$\mathcal{L}_{CE}(\hat{y}, y)$$
$$-\nabla_{\theta} \mathcal{L}_{CE}(\hat{y}, y)$$

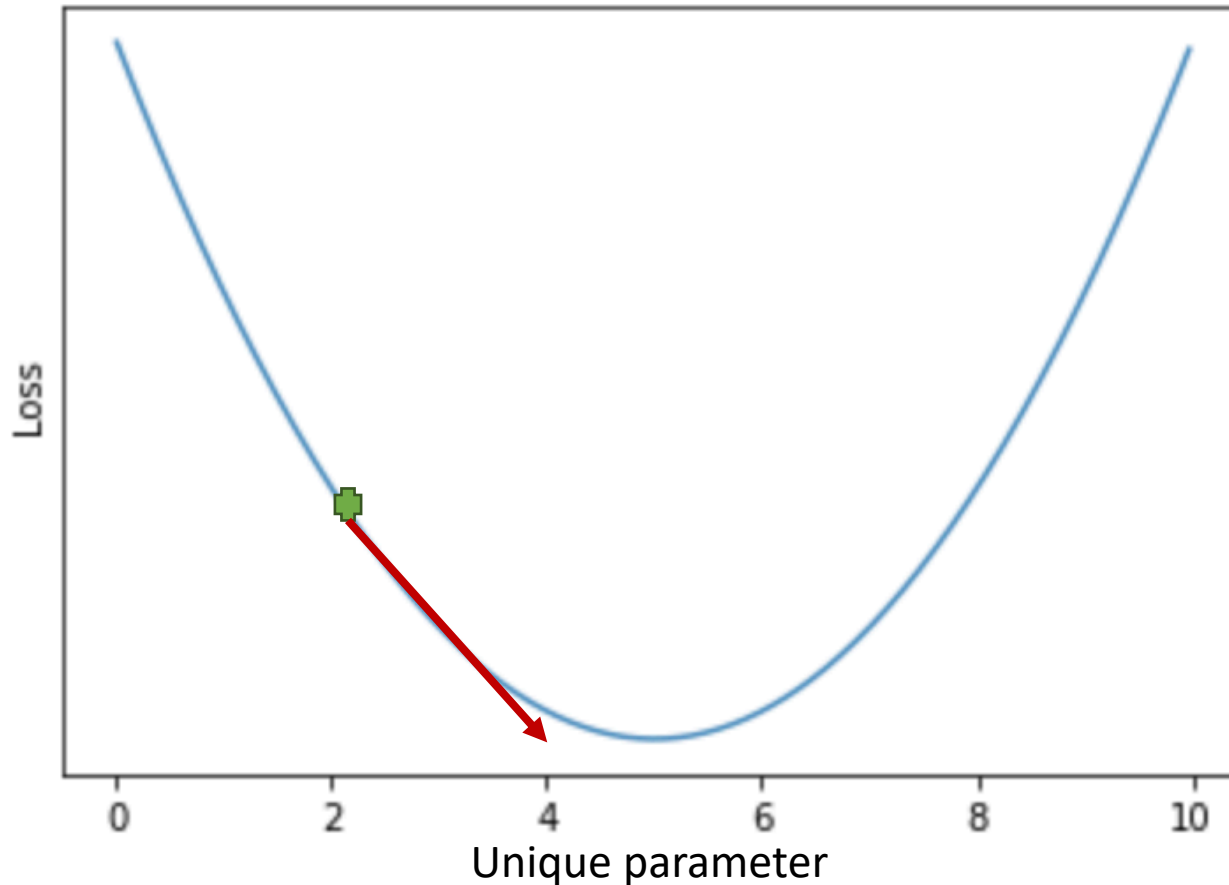


Minimization of a function



Optimize the network to minimize the loss with respect to all neurons θ :

$$\mathcal{L}_{CE}(\hat{y}, y)$$
$$-\nabla_{\theta} \mathcal{L}_{CE}(\hat{y}, y)$$

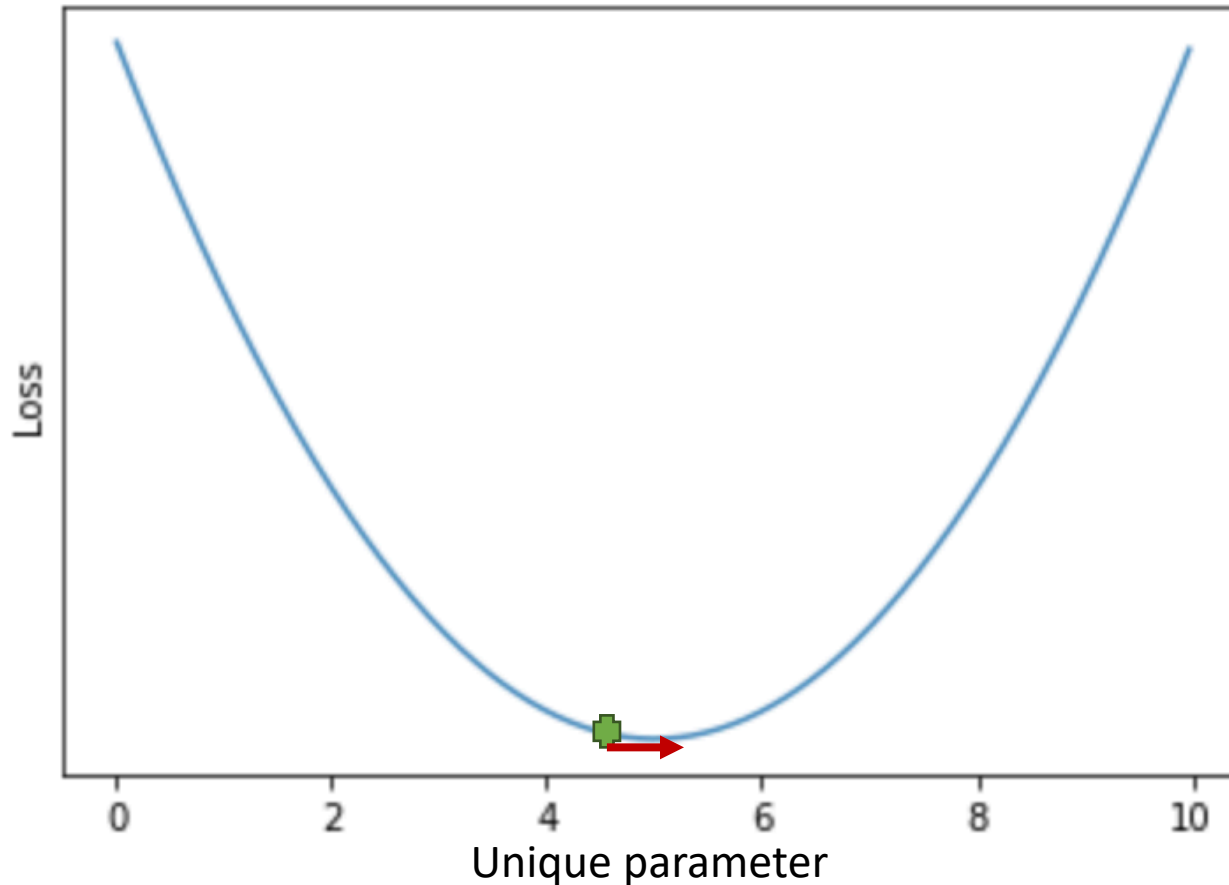


Minimization of a function



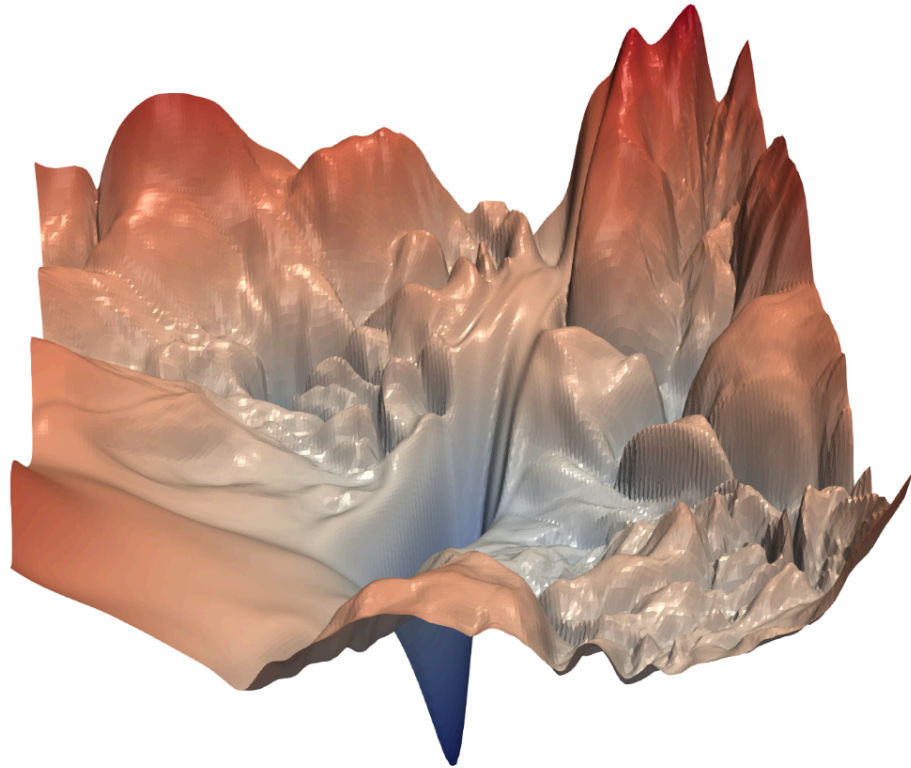
Optimize the network to minimize the loss with respect to all neurons θ :

$$\mathcal{L}_{CE}(\hat{y}, y)$$
$$-\nabla_{\theta} \mathcal{L}_{CE}(\hat{y}, y)$$





Million of parameters to optimize together!
Highly non-convex!

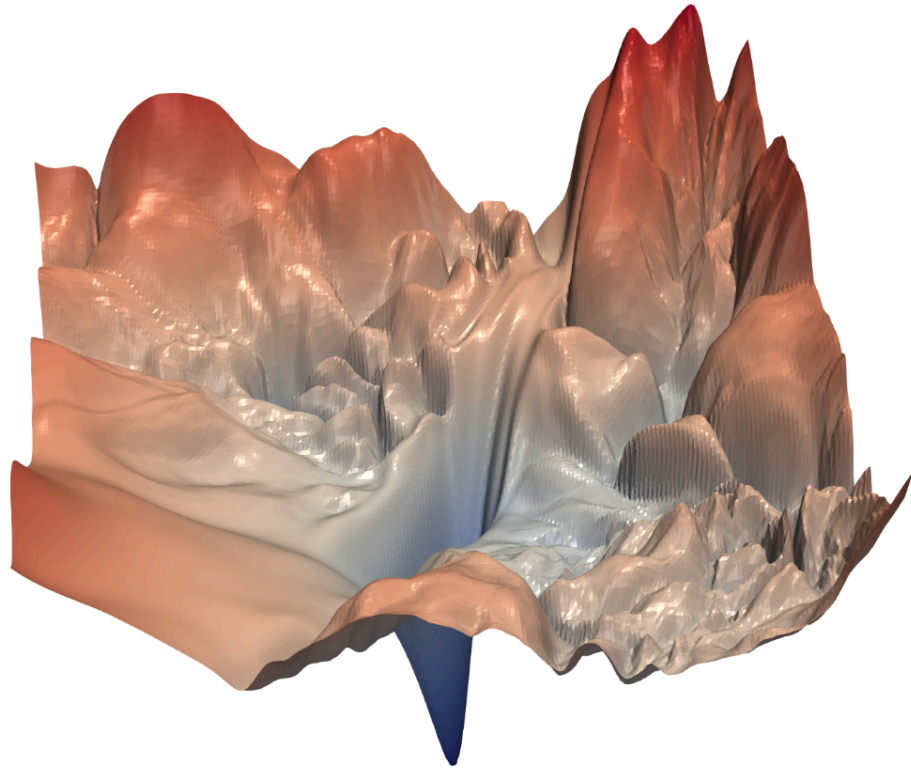


Multiple dimensional derivatives are called **gradients**

Minimization of a function



In theory a lot of bad local minima, but in practice a DNN usually find a good local minima close to optimum [[LeCun et al. Nature 2015](#)].

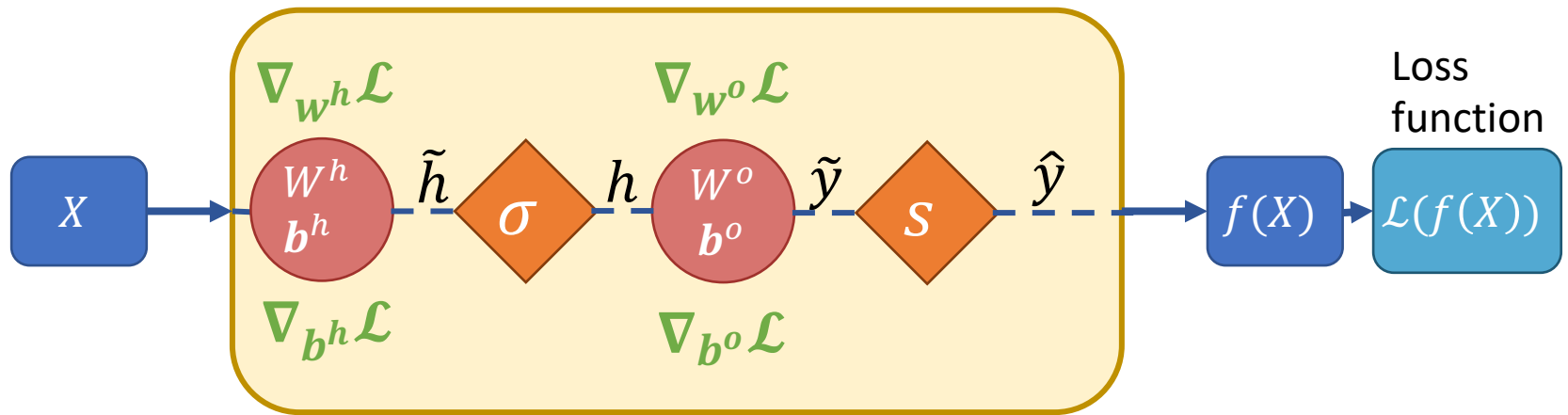


Although, this is debated [[Goldblum et al. ICLR 2020](#)]

Needed Gradients



We need all parameters gradients in **green**

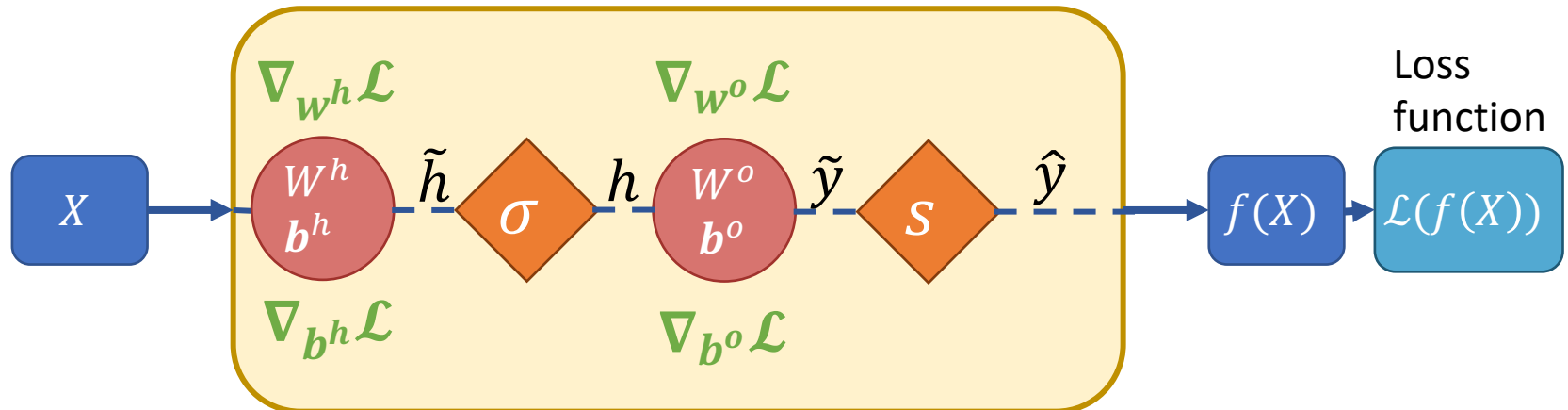


Stochastic Gradient Descent



1. Initialize randomly the parameters θ
2. For each epoch do
 1. Select a random sample of the data
 2. Forward
3. Compute gradients $\nabla_{\theta} \mathcal{L}$ for each parameter θ
4. Update parameters $\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}$

η is the **learning rate**.





$$f \circ g(a) = f(b) = c$$

Given scalars:

$$\frac{\partial c}{\partial a} = \frac{\partial c}{\partial b} \frac{\partial b}{\partial a}$$

Given vectors,
element-wise:

$$\frac{\partial c_i}{\partial a_k} = \sum_j \frac{\partial c_i}{\partial b_j} \frac{\partial b_j}{\partial a_k}$$

Given vectors,
vector-wise:

$$\nabla_a c = \nabla_b c \nabla_a b^T$$

In [denominator layout](#) (\neq numerator layout):

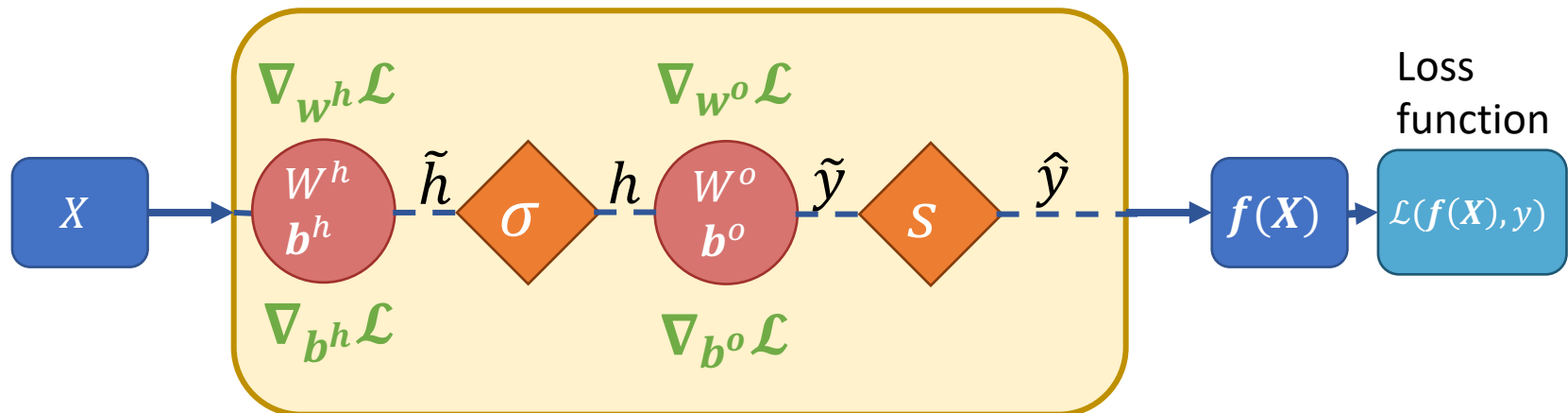
$$\mathbf{a} = \begin{bmatrix} a_1 \\ \vdots \\ a_{n_a} \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_{n_b} \end{bmatrix} \quad \mathbf{c} = \begin{bmatrix} c_1 \\ \vdots \\ c_{n_c} \end{bmatrix}$$

$$\nabla_b c = \begin{bmatrix} \frac{\partial c_1}{\partial b_1} & \dots & \frac{\partial c_{n_c}}{\partial b_1} \\ \vdots & & \vdots \\ \frac{\partial c_1}{\partial b_{n_b}} & \dots & \frac{\partial c_{n_c}}{\partial b_{n_b}} \end{bmatrix}$$



Partial derivative of the cross-entropy loss with respect to (w.r.t.) the probabilities:

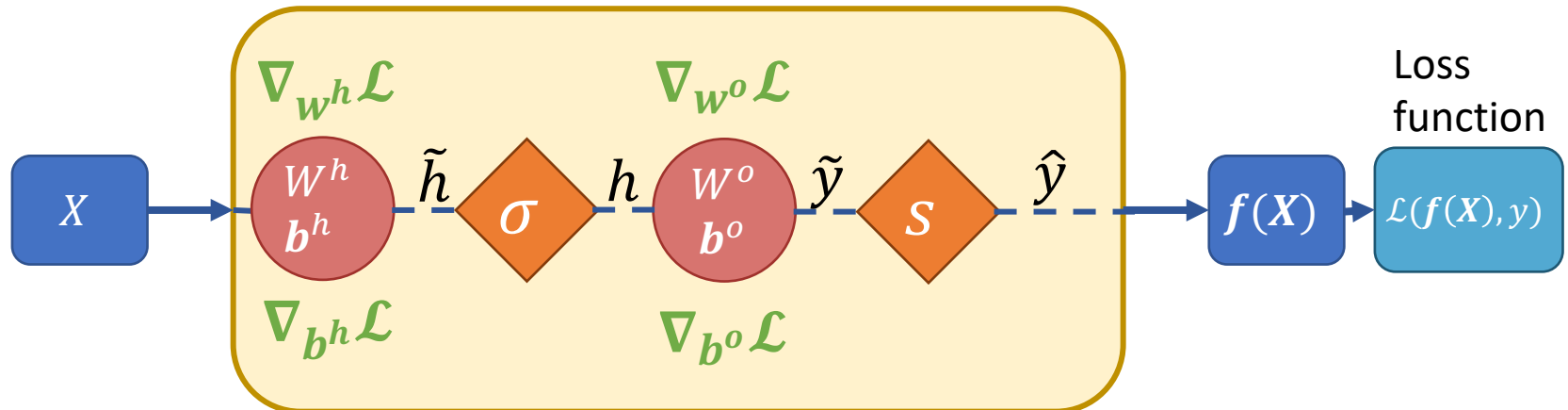
$$\frac{\partial \mathcal{L}(f(x), y)}{\partial f(x)_i} = \frac{\partial -\log f(x)_y}{\partial f(x)_i} = \frac{-1_{y=i}}{f(x)_y}, \text{ for simplicity } \frac{\partial \mathcal{L}}{\partial f(x)_i}$$



Backpropagation



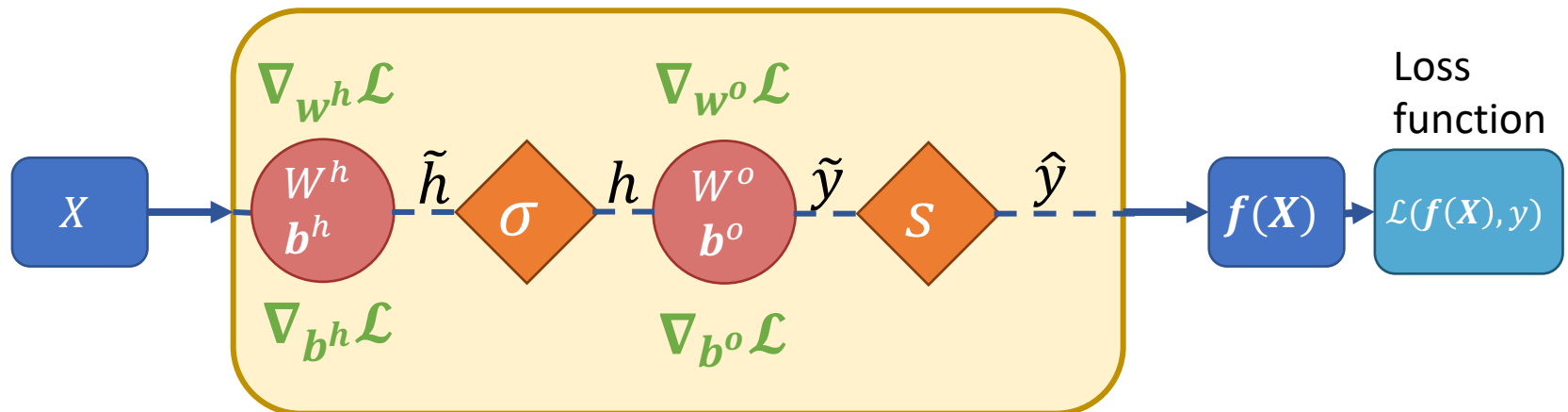
$$\frac{\partial \mathcal{L}}{\partial \tilde{y}_i} = \sum_j \frac{\partial \mathcal{L}}{\partial f(x)_j} \frac{\partial f(x)_j}{\partial \tilde{y}_i}$$



Backpropagation



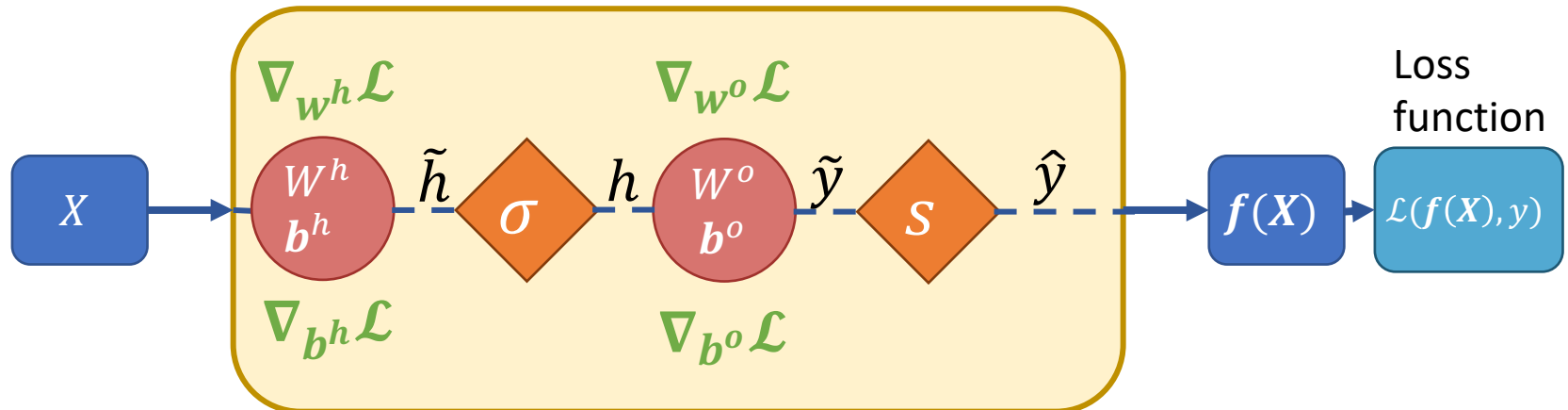
$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \tilde{y}_i} &= \sum_j \frac{\partial \mathcal{L}}{\partial f(\mathbf{x})_j} \frac{\partial f(\mathbf{x})_j}{\partial \tilde{y}_i} \\ &= \sum_j \frac{-1_{y=j}}{\partial f(\mathbf{x})_y} \frac{\partial \text{softmax}(\tilde{\mathbf{y}})_j}{\partial \tilde{y}_i}\end{aligned}$$



Backpropagation



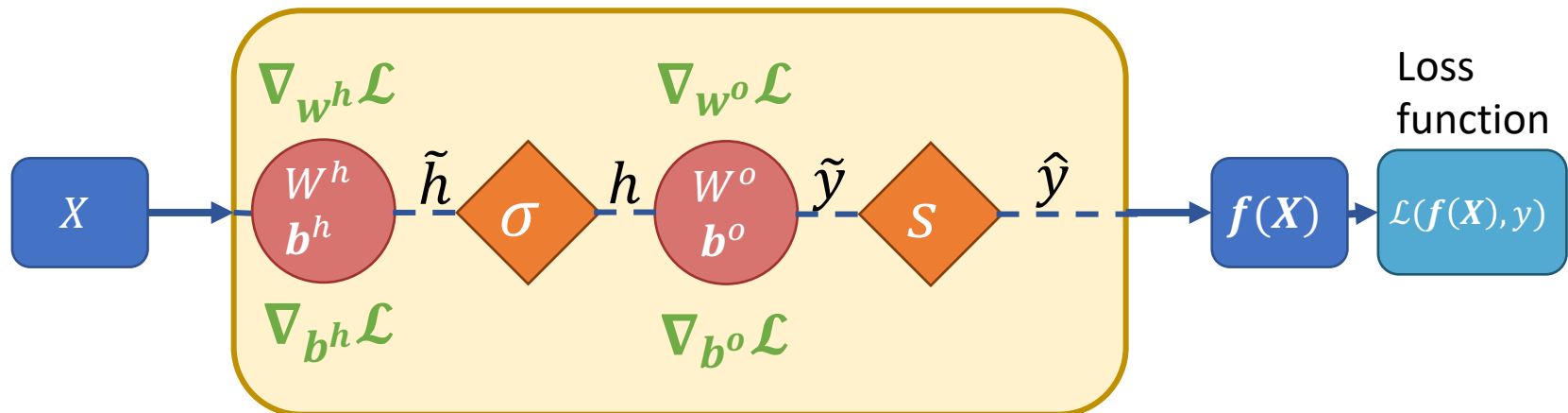
$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \tilde{y}_i} &= \sum_j \frac{\partial \mathcal{L}}{\partial f(\mathbf{x})_j} \frac{\partial f(\mathbf{x})_j}{\partial \tilde{y}_i} \\ &= \sum_j \frac{-1_{y=j}}{\partial f(\mathbf{x})_y} \frac{\partial \text{softmax}(\tilde{\mathbf{y}})_j}{\partial \tilde{y}_i} \\ &= \begin{cases} \frac{-1}{f(\mathbf{x})_y} \text{softmax}(\tilde{\mathbf{y}})_y (1 - \text{softmax}(\tilde{\mathbf{y}})_y) & \text{if } i = y \\ \frac{1}{f(\mathbf{x})_y} \text{softmax}(\tilde{\mathbf{y}})_y \text{softmax}(\tilde{\mathbf{y}})_i & \text{if } i \neq y \end{cases}\end{aligned}$$



Backpropagation



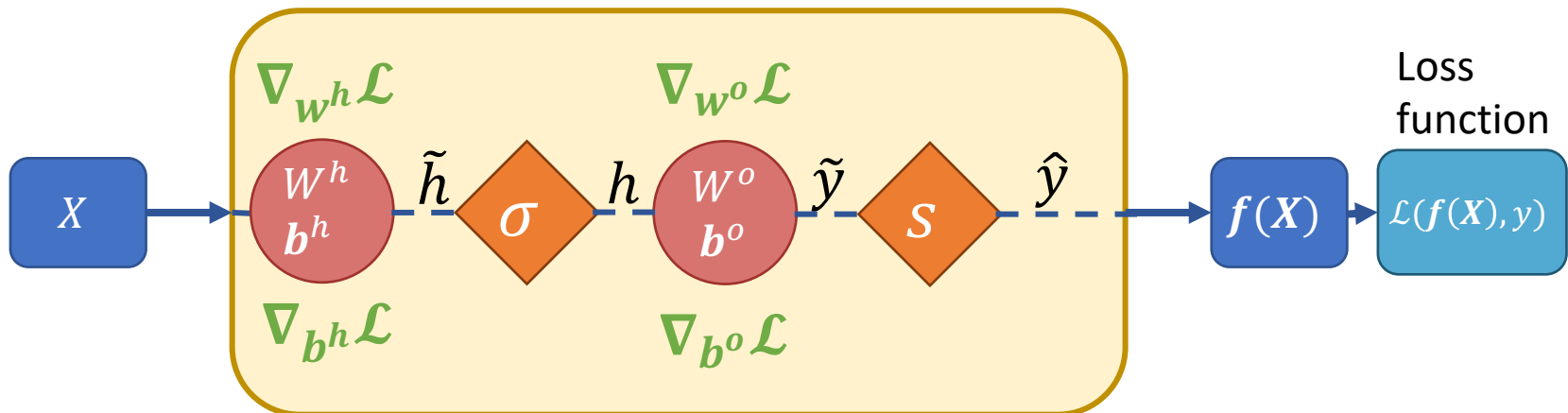
$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \tilde{y}_i} &= \sum_j \frac{\partial \mathcal{L}}{\partial \mathbf{f}(\mathbf{x})_j} \frac{\partial \mathbf{f}(\mathbf{x})_j}{\partial \tilde{y}_i} \\ &= \sum_j \frac{-1_{y=j}}{\mathbf{f}(\mathbf{x})_y} \frac{\partial \text{softmax}(\tilde{\mathbf{y}})_j}{\partial \tilde{y}_i} \\ &= \begin{cases} \frac{-1}{\mathbf{f}(\mathbf{x})_y} \text{softmax}(\tilde{\mathbf{y}})_y (1 - \text{softmax}(\tilde{\mathbf{y}})_y) & \text{if } i = y \\ \frac{1}{\mathbf{f}(\mathbf{x})_y} \text{softmax}(\tilde{\mathbf{y}})_y \text{softmax}(\tilde{\mathbf{y}})_i & \text{if } i \neq y \end{cases} \\ &= \begin{cases} -1 + \mathbf{f}(\mathbf{x})_y & \text{if } i = y \\ \mathbf{f}(\mathbf{x})_i & \text{if } i \neq y \end{cases}\end{aligned}$$



Backpropagation



$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \tilde{y}_i} &= \sum_j \frac{\partial \mathcal{L}}{\partial f(\mathbf{x})_j} \frac{\partial f(\mathbf{x})_j}{\partial \tilde{y}_i} \\ &= \sum_j \frac{-1_{y=j}}{\partial f(\mathbf{x})_y} \frac{\partial \text{softmax}(\tilde{\mathbf{y}})_j}{\partial \tilde{y}_i} \\ &= \begin{cases} \frac{-1}{f(\mathbf{x})_y} \text{softmax}(\tilde{\mathbf{y}})_y (1 - \text{softmax}(\tilde{\mathbf{y}})_y) & \text{if } i = y \\ \frac{1}{f(\mathbf{x})_y} \text{softmax}(\tilde{\mathbf{y}})_y \text{softmax}(\tilde{\mathbf{y}})_i & \text{if } i \neq y \end{cases} \\ &= \begin{cases} -1 + f(\mathbf{x})_y & \text{if } i = y \\ f(\mathbf{x})_i & \text{if } i \neq y \end{cases} \end{aligned} \quad \nabla_{\tilde{\mathbf{y}}} \mathcal{L} = \mathbf{f}(\mathbf{x}) - \mathbf{e}(y) \text{ with one-hot encoding of the target}$$



Backpropagation

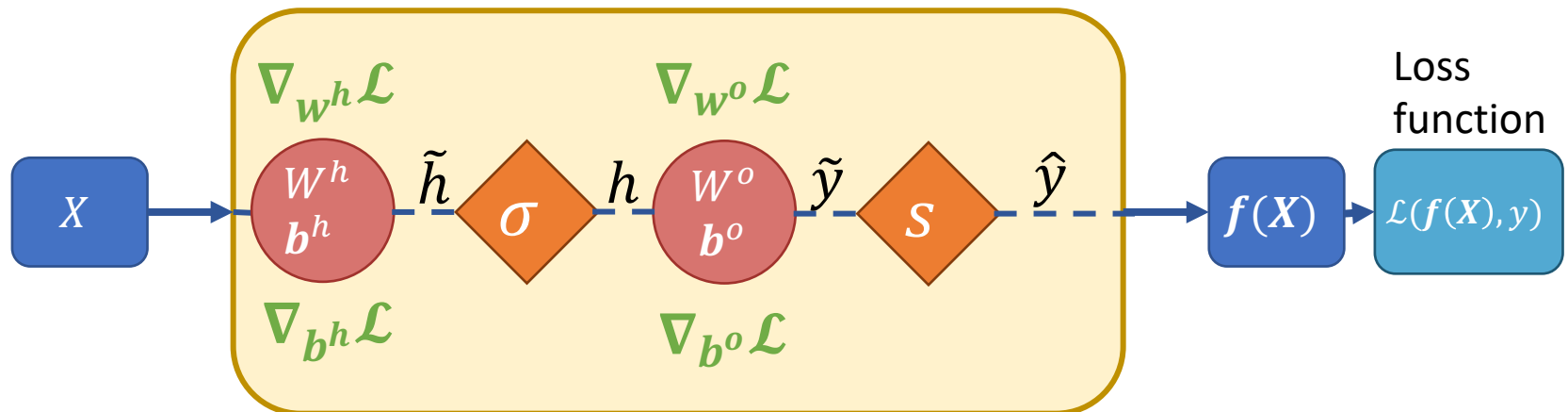


$$\nabla_{\tilde{y}} \mathcal{L} = \mathbf{f}(x) - \mathbf{e}(y)$$

$$\nabla_{b^o} \mathcal{L} = \nabla_{\tilde{y}} \mathcal{L}$$

$$\nabla_{w^o} \mathcal{L} = \nabla_{\tilde{y}} \mathcal{L} \cdot \mathbf{h}^T$$

$$\tilde{\mathbf{y}} = \mathbf{W}^o \mathbf{h} + \mathbf{b}^o$$



Backpropagation



$$\nabla_{\tilde{y}} \mathcal{L} = \mathbf{f}(\mathbf{x}) - \mathbf{e}(y)$$

$$\nabla_{b^o} \mathcal{L} = \nabla_{\tilde{y}} \mathcal{L}$$

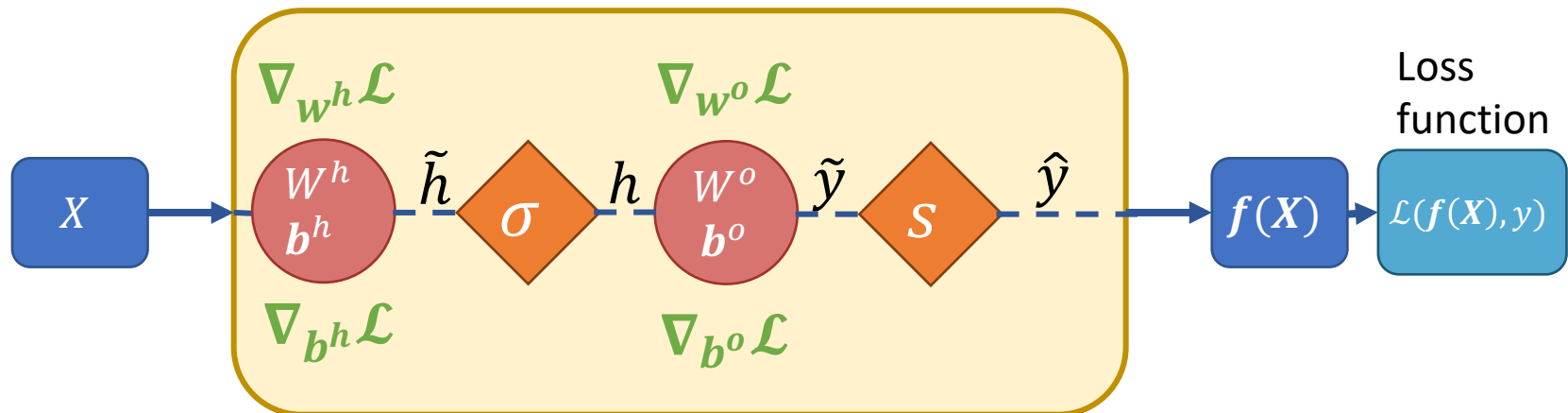
$$\nabla_{w^o} \mathcal{L} = \nabla_{\tilde{y}} \mathcal{L} \cdot \mathbf{h}^T$$

$$\nabla_{\mathbf{h}} \mathcal{L} = \mathbf{W}^{oT} \cdot \nabla_{\tilde{y}} \mathcal{L}$$

$$\nabla_{\tilde{\mathbf{h}}} \mathcal{L} = \nabla_{\mathbf{h}} \mathcal{L} \odot \sigma'(\tilde{\mathbf{h}})$$

\odot is the element-wise multiplication (Hadamard product).

Sigmoid is applied element-wise, thus the gradient also.



Backpropagation



$$\nabla_{\tilde{y}} \mathcal{L} = \mathbf{f}(x) - \mathbf{e}(y)$$

$$\nabla_{b^o} \mathcal{L} = \nabla_{\tilde{y}} \mathcal{L}$$

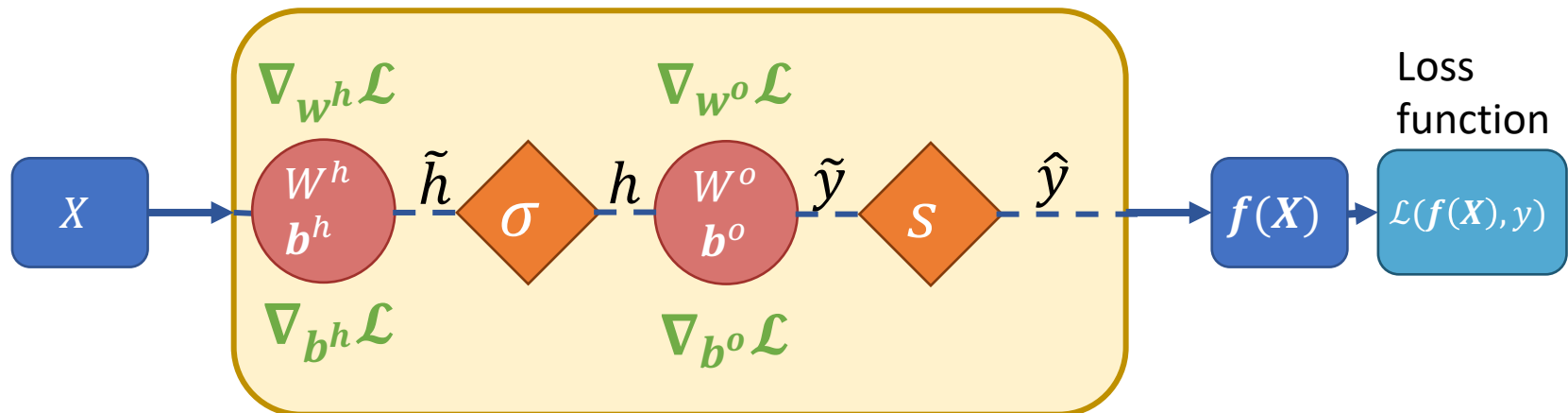
$$\nabla_{w^o} \mathcal{L} = \nabla_{\tilde{y}} \mathcal{L} \cdot \mathbf{h}^T$$

$$\nabla_{\mathbf{h}} \mathcal{L} = \mathbf{W}^{oT} \cdot \nabla_{\tilde{y}} \mathcal{L}$$

$$\nabla_{\tilde{\mathbf{h}}} \mathcal{L} = \nabla_{\mathbf{h}} \mathcal{L} \odot \sigma'(\tilde{\mathbf{h}})$$

$$\nabla_{b^h} \mathcal{L} = ?$$

$$\nabla_{W^h} \mathcal{L} = ?$$



Backpropagation



$$\nabla_{\tilde{y}} \mathcal{L} = \mathbf{f}(\mathbf{x}) - \mathbf{e}(y)$$

$$\nabla_{b^o} \mathcal{L} = \nabla_{\tilde{y}} \mathcal{L}$$

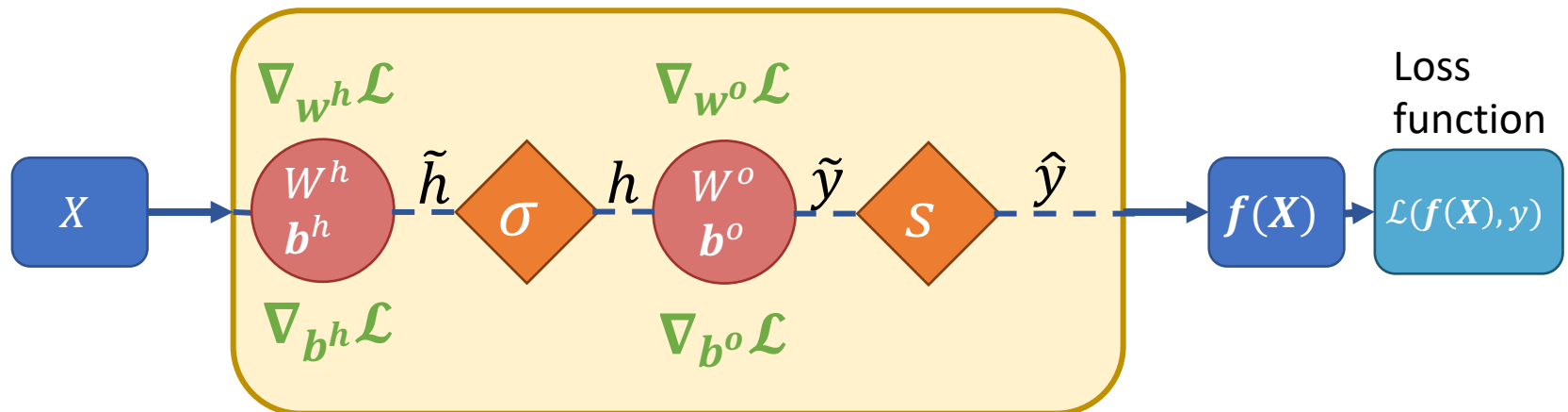
$$\nabla_{w^o} \mathcal{L} = \nabla_{\tilde{y}} \mathcal{L} \cdot \mathbf{h}^T$$

$$\nabla_{\mathbf{h}} \mathcal{L} = \mathbf{W}^{oT} \cdot \nabla_{\tilde{y}} \mathcal{L}$$

$$\nabla_{\tilde{\mathbf{h}}} \mathcal{L} = \nabla_{\mathbf{h}} \mathcal{L} \odot \sigma'(\tilde{\mathbf{h}})$$

$$\nabla_{b^h} \mathcal{L} = \nabla_{\tilde{\mathbf{h}}} \mathcal{L}$$

$$\nabla_{W^h} \mathcal{L} = \nabla_{\tilde{\mathbf{h}}} \mathcal{L} \cdot \mathbf{x}^T$$



Backpropagation



$$\nabla_{\tilde{y}} \mathcal{L} = \mathbf{f}(\mathbf{x}) - \mathbf{e}(y)$$

$$\nabla_{b^o} \mathcal{L} = \nabla_{\tilde{y}} \mathcal{L}$$

$$\nabla_{w^o} \mathcal{L} = \nabla_{\tilde{y}} \mathcal{L} \cdot \mathbf{h}^T$$

$$\nabla_{\mathbf{h}} \mathcal{L} = \mathbf{W}^{oT} \cdot \nabla_{\tilde{y}} \mathcal{L}$$

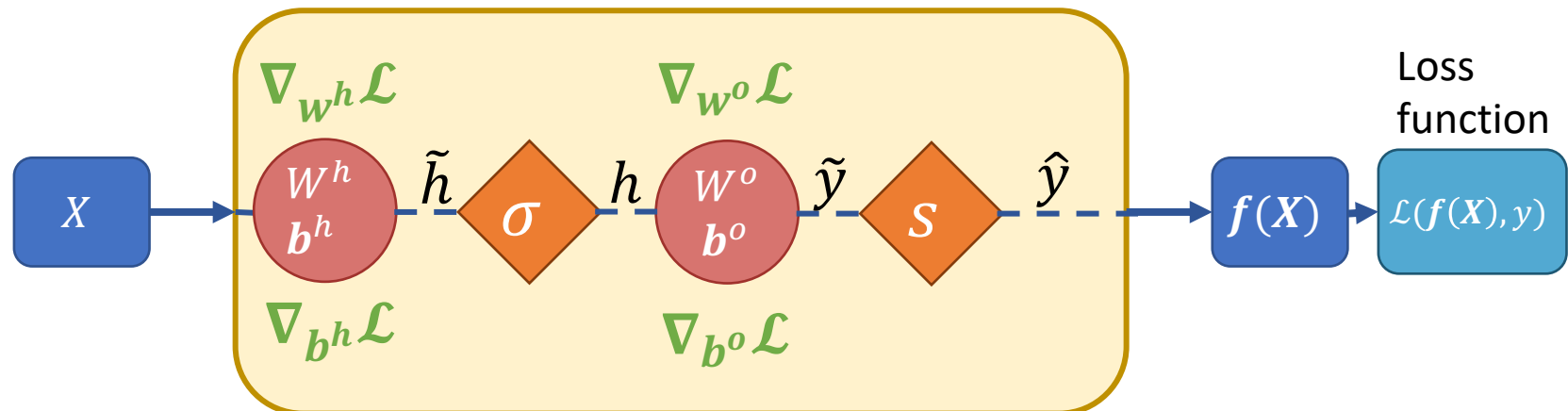
$$\nabla_{\tilde{\mathbf{h}}} \mathcal{L} = \nabla_{\mathbf{h}} \mathcal{L} \odot \sigma'(\tilde{\mathbf{h}})$$

$$\nabla_{b^h} \mathcal{L} = \nabla_{\tilde{\mathbf{h}}} \mathcal{L}$$

$$\nabla_{W^h} \mathcal{L} = \nabla_{\tilde{\mathbf{h}}} \mathcal{L} \cdot \mathbf{x}^T$$

When in doubt, look at the shape.

$\nabla_W \mathcal{L}$ must have the shape of W because of the update rule $W \leftarrow W - \nabla_W \mathcal{L}$



Backpropagation



$$\nabla_{\tilde{y}} \mathcal{L} = \mathbf{f}(\mathbf{x}) - \mathbf{e}(y)$$

$$\nabla_{b^o} \mathcal{L} = \nabla_{\tilde{y}} \mathcal{L}$$

$$\nabla_{w^o} \mathcal{L} = \nabla_{\tilde{y}} \mathcal{L} \cdot \mathbf{h}^T$$

$$\nabla_{\mathbf{h}} \mathcal{L} = \mathbf{W}^{oT} \cdot \nabla_{\tilde{y}} \mathcal{L}$$

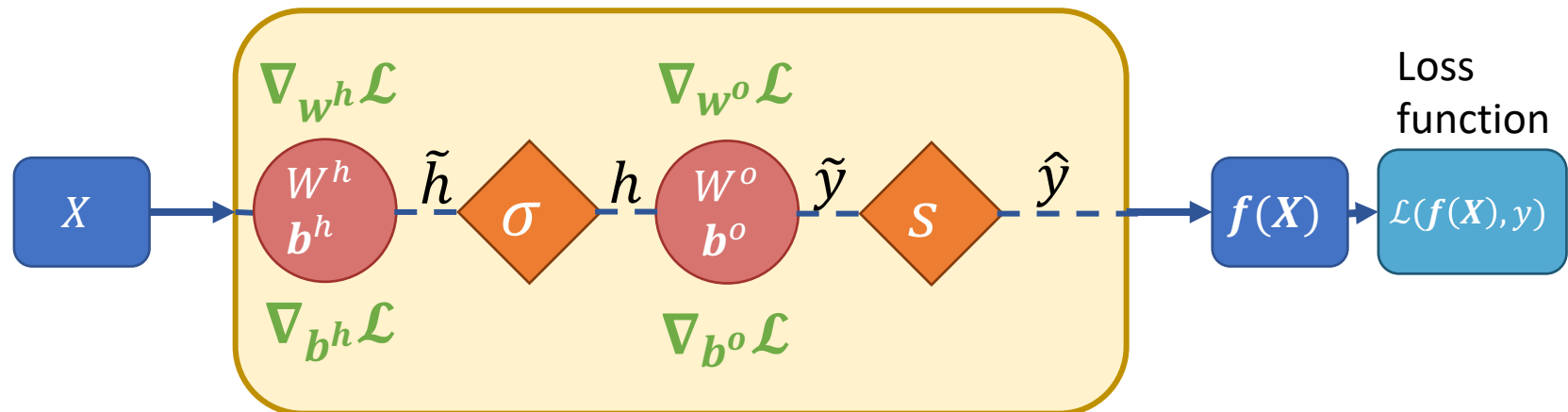
$$\nabla_{\tilde{\mathbf{h}}} \mathcal{L} = \nabla_{\mathbf{h}} \mathcal{L} \odot \sigma'(\tilde{\mathbf{h}})$$

$$\nabla_{b^h} \mathcal{L} = \nabla_{\tilde{\mathbf{h}}} \mathcal{L}$$

$$\nabla_{W^h} \mathcal{L} = \nabla_{\tilde{\mathbf{h}}} \mathcal{L} \cdot \mathbf{x}^T$$

To have huge speed-up:

1. Re-use previous gradients
2. Save tensors during forward



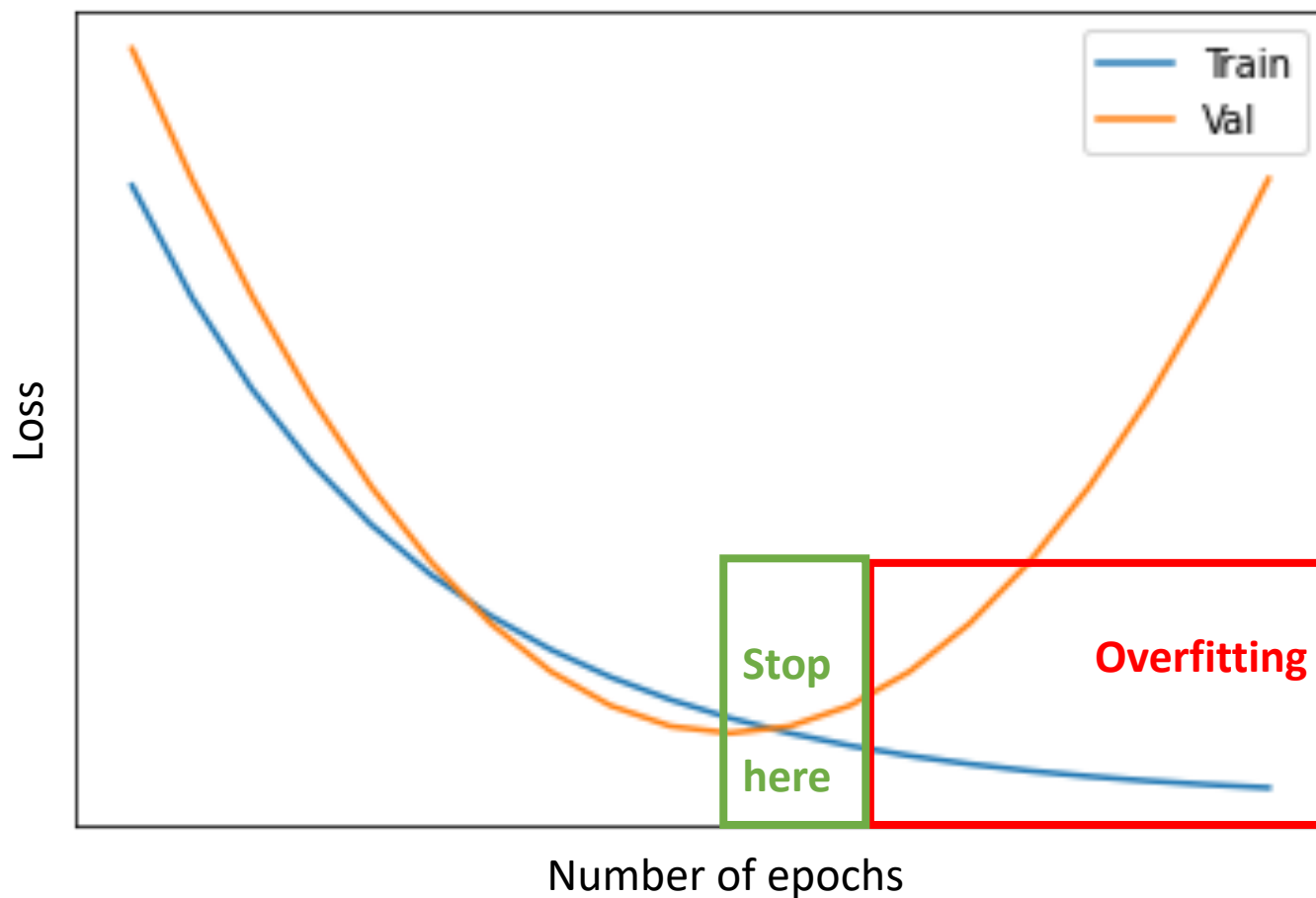
Tips & Tricks

When to Stop?



Split data in **train / val / test**

Stop when a criterion (loss, accuracy, f1, etc.) stop improving on **validation set**





Batch Gradient Descent: one forward & backward on the whole dataset

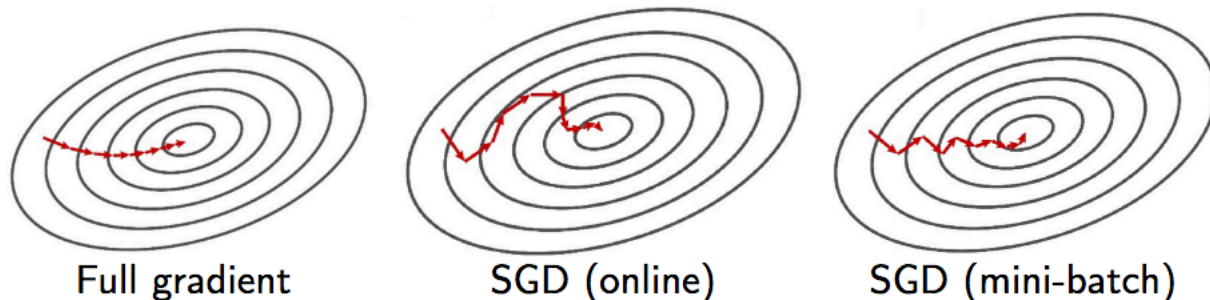
- Better gradient estimation
- GPU parallelism
- Impracticable to fit large dataset in VRAM

Stochastic Gradient Descent: one forward & backward per sample

- Easy to fit in VRAM
- Add noise that may improve generalization
- Add too much noise
- Slow

Mini-Batch Gradient Descent: one forward & backward per group of samples

- Trade-off between both
- Learning rate should be proportional to batch size, e.g. batch size 32- \rightarrow 64, lr 0.1- \rightarrow 0.2





A fully connected layer in forward pass was:

$$\mathbf{h} = \mathbf{W}\mathbf{x} + \mathbf{b}$$

Considering that $\mathbf{x} \in \mathbb{R}^n$ was the features of a unique sample.

But with a batch size $\mathbf{X} \in \mathbb{R}^{b \times n}$, thus:

$$\mathbf{h} = \mathbf{X}\mathbf{W}^T + \mathbf{b}$$

With \mathbf{W} defined as before.



$$\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}$$

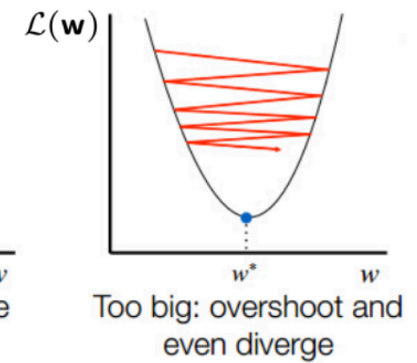
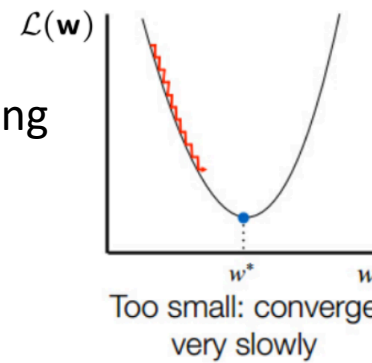
Controls the rate of change.

Too high:

→ Cannot converge, but diverge, reduce overfitting

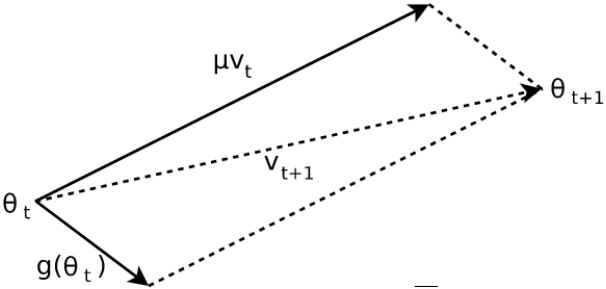
Too low:

→ Super slow, stuck in bad local minima



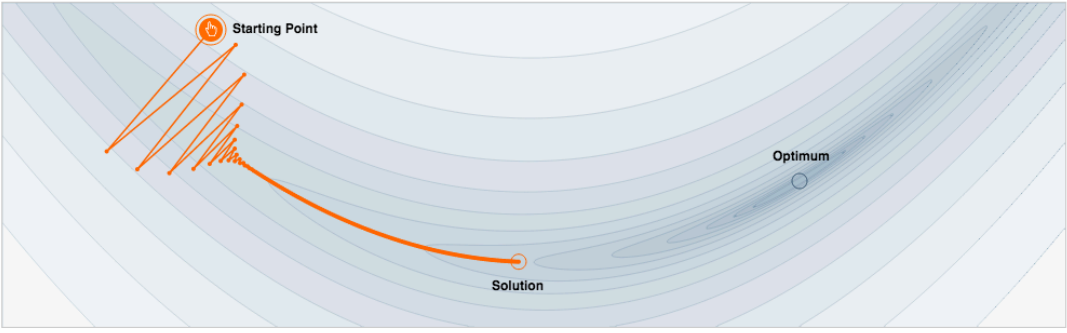
Tip: Start with high learning rate, and decreases it through time

SGD with Momentum

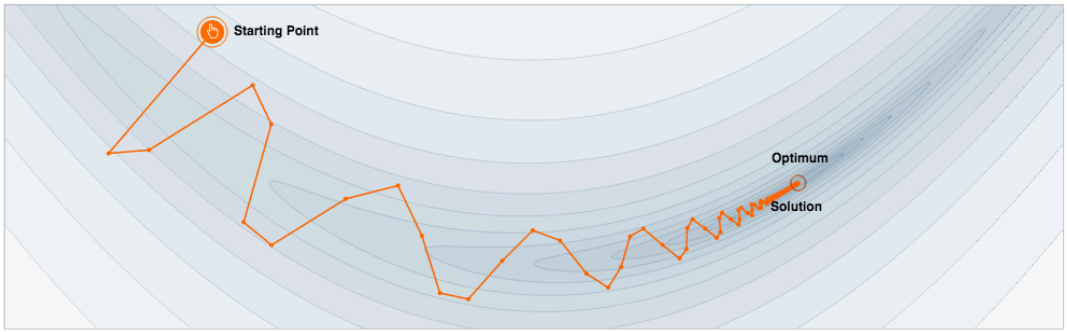


$$v \leftarrow \alpha v - (1 - \alpha) \nabla_{\theta} \mathcal{L}$$

$$\theta \leftarrow \theta + v$$



We often think of Momentum as a means of dampening oscillations and speeding up the iterations, leading to faster convergence. But it has other interesting behavior. It allows a larger range of step-sizes to be used, and creates its own oscillations. What is going on?



We often think of Momentum as a means of dampening oscillations and speeding up the iterations, leading to faster convergence. But it has other interesting behavior. It allows a larger range of step-sizes to be used, and creates its own oscillations. What is going on?

[Why Momentum Really works, on distill.pub](#)

Optimizers

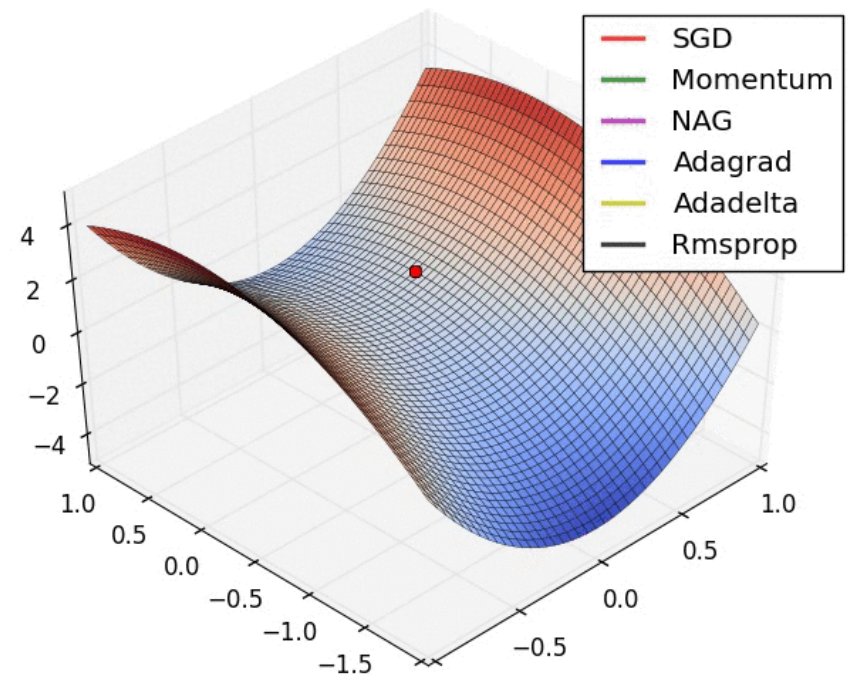
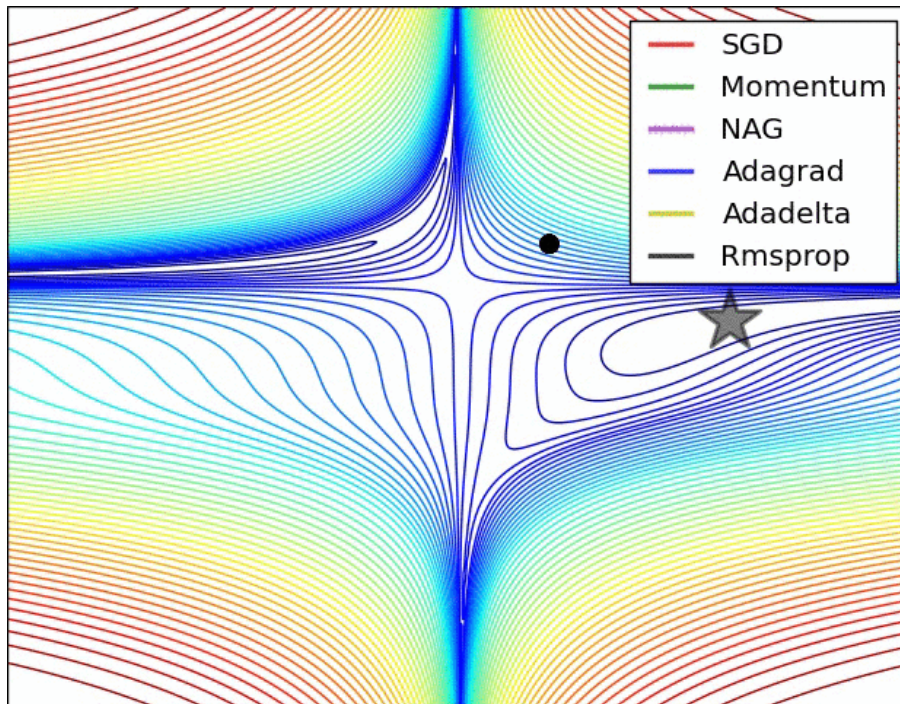


Modern optimizers have an **adaptive learning rate** per parameter based on gradient statistics.

→ Especially useful on saddle point

→ The most famous is Adam

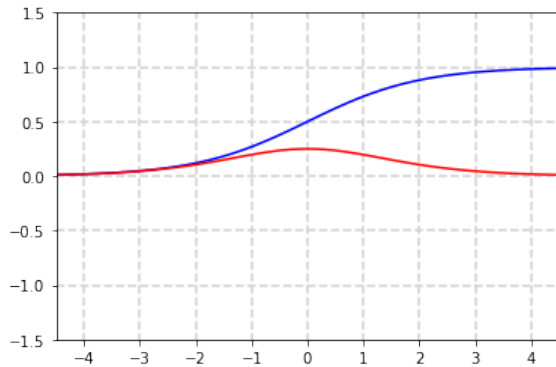
But a well-tuned SGD with momentum can be the best [[Wilson et al. 2017](#)].



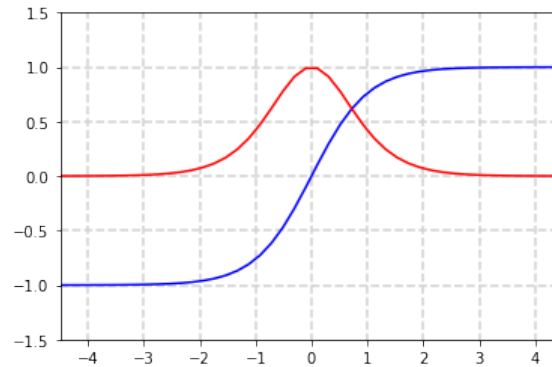
Great overview of gradient-descent based optimizers by [Ruder](#).



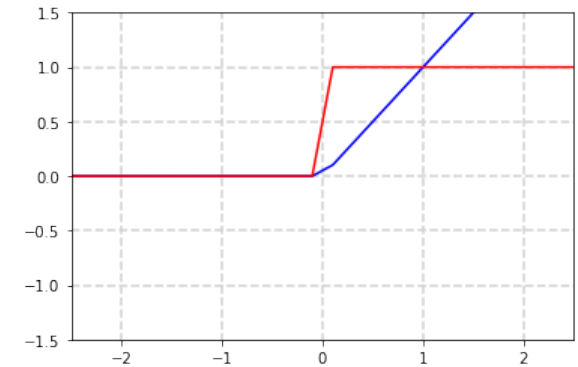
Function and their derivative



Sigmoid



Tanh



ReLU

Sigmoid and tanh **saturates** at small and large values.

→ Gradient is zero, no learning

→ Avoid these old-school activations

ReLU is zero if $x \leq 0$:

→ **Dying neurons** with zero-output and thus zero-gradient

→ If it happens, use a **Leaky ReLU** $LReLU(x) = \begin{cases} x & \text{if } x > 0 \\ \epsilon x & \text{otherwise} \end{cases}$



Initializing the biases \mathbf{b}^h and \mathbf{b}^o to very small values.

→ Helpful to avoid dying neurons with ReLU

Initializing the weights \mathbf{W}^h and \mathbf{W}^o to:

- **Zero-weights**

- No learning because gradient w.r.t input is also zero

- **Constant weights**

- Symmetry where two hidden neurons are connected to the same inputs, they learn the same pattern!

- **Large values**

- Risk gradient explosion

- **He / Glorot initialization**

- Normalize weights to avoid explosion with large number of outgoing connections



$$\mathcal{L} = \mathcal{L}_1(x, y) + \lambda \mathcal{L}_2(x, y) + \beta \mathcal{L}(\theta)$$

Auxiliary loss

You can combine multiple losses, each has been batch averaged.

Main loss (usually the classification loss) has a factor of 1.

Auxiliary losses have a factor as hyperparameter

- Need to find optimal through cross-validation
- Should often ensure that losses are more or less in the same values

Regularization losses usually only takes the parameters as argument

- Most common is **weight decay** $\beta \sum_i \|\theta_i\|^2$
- Often useful to put a prior on the parameters

Small break,
then coding session!