# DEEP NEURAL NETWORK

## Deep Learning for Computer Vision

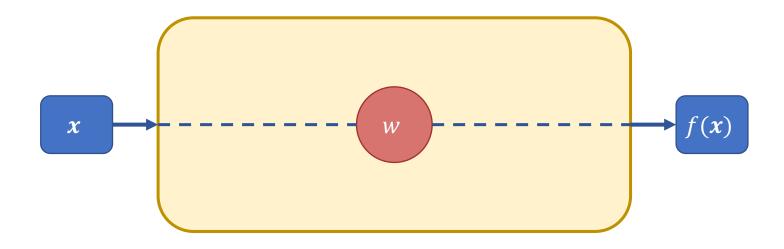Arthur Douillard

# Deep Neural Networks

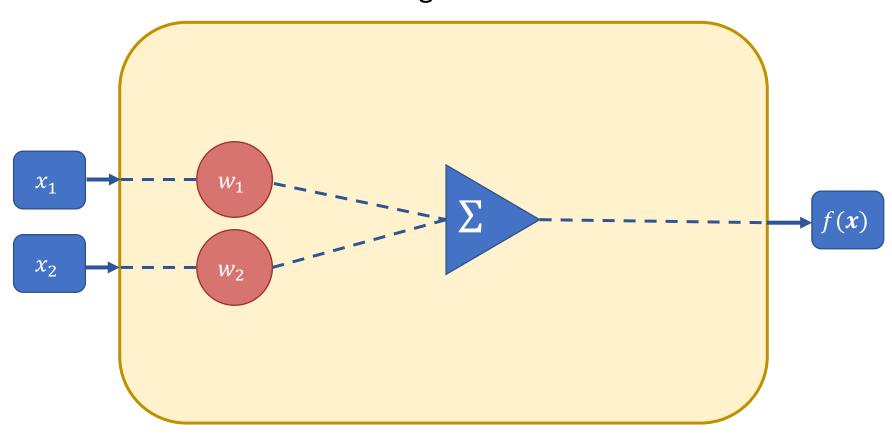A neural network <u>approximates</u> a function

Linear regression



Useful to predict continuous variables
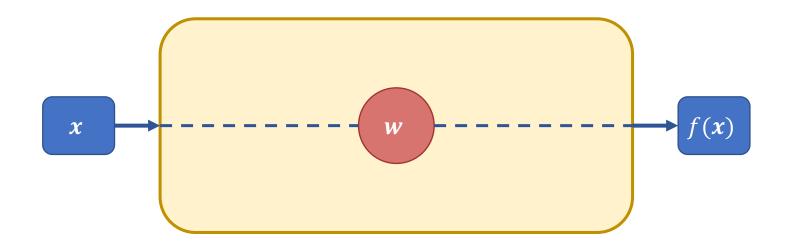
$$f(x) = wx$$

Linear regression

The "connection" often depicted are only multiplications and additions
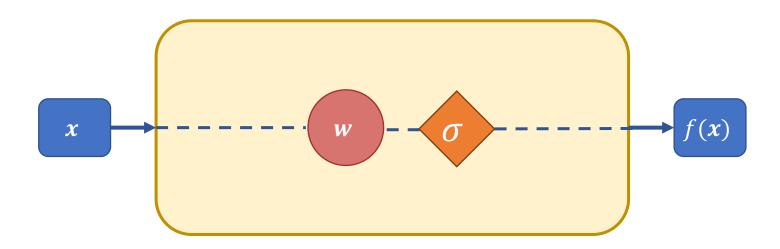
$$f(x) = w_1 x_1 + w_2 x_2$$

Linear regression

$$f(x) = wx$$

## Single Neuron

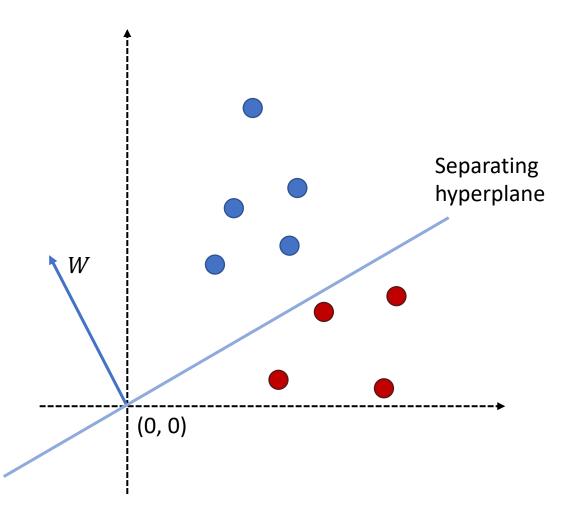

**Non-linear activation** to determine how much a neuron "fires"

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

Optimize so that $W$
is orthogonal to
the separating hyperplane

$W$

Separating
hyperplane

(0, 0)

Need bias!

*W*

(0, 0)

Separating
hyperplane

Single Neuron



With a **bias**

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

Separating hyperplane

$W$

$b$

(0, 0)

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



(0, 0)

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**

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

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

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

# Multi-Layer Perceptron



$$\widetilde{\boldsymbol{h}} = \boldsymbol{W}^h \boldsymbol{x} + \boldsymbol{b}^h$$

Embeddings/features



$$\boldsymbol{h} = \sigma(\widetilde{\boldsymbol{h}})$$

# Multi-Layer Perceptron

Logits!



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

# Multi-Layer Perceptron



$$\widehat{\boldsymbol{y}} = \sigma(\widetilde{\boldsymbol{y}})$$

Considering hidden dimension of 3.

Hidden layer can define 3 hyperplanes

(0, 0)

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

## Multi-Layer Perceptron



$$\widetilde{h} = W^h x + b^h$$
$$h = \sigma(\widetilde{h})$$
$$\widetilde{y} = W^o h + b^o$$
$$\widehat{y} = \sigma(\widetilde{y})$$

## Multi-Layer Perceptron



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

$$h = \sigma(\widetilde{h}) \qquad \rightarrow \text{features / embeddings}$$

$$\widetilde{y} = W^o h + b^o \qquad \rightarrow \text{logits}$$

$$\hat{y} = \sigma(\widetilde{y}) \qquad \rightarrow \text{Model predictions}$$

**Function** and their **derivative**



| Sigmoid | Hyperbolic tangent | Rectified Linear Unit |
|---------|--------------------|-----------------------|

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

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

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

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

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

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

**Function** and their **derivative**



| Sigmoid | Hyperbolic tangent | Rectified Linear Unit |

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

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

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

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

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

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

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

→ Composition of affine functions is an affine function

# Output activations

Binary classification: **sigmoid**
→ Applied element-wise
→ Range [0, 1]

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

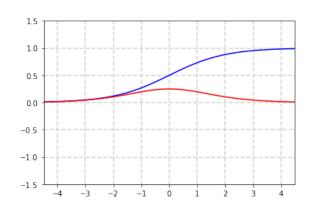$$\frac{d}{dx}\sigma(x) = \sigma(x)(1 - \sigma(x))$$

Multi-Class classification: **softmax**
→ Applied over a vector
→ Range [0, 1] per element
→ Sum to 1 for the vector
    → Probability distribution

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

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

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

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 pass

$X$ $\quad$ $W^h$, $\boldsymbol{b}^h$ $\quad$ $\sigma$ $\quad$ $W^o$, $\boldsymbol{b}^o$ $\quad$ $S$ $\quad$ $f(X)$ $\quad$ $\mathcal{L}(f(X))$

Loss function

Backward pass

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: $-\log 0.8 = 0.2231 \dots$

# One-Hot

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:

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

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

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

Optimize the network to minimize the loss with respect to all neurons $\theta$:

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

Optimize the network to minimize the loss with respect to all neurons $\theta$ :

$$\mathcal{L}_{CE}(\widehat{y}, \; y)$$

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

Optimize the network to minimize the loss with respect to all neurons $\theta$:

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

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



Multiple dimensional derivatives are called **gradients**

Visualizing the Loss Landscape of Neural Nets, [Li et al. NeurIPS2018]

We need all parameters gradients in **green**

# Stochastic Gradient Descent

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

$\eta$ 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:
$$\frac{\partial c_i}{\partial a_k} = \sum_i \frac{\partial c_i}{\partial b_j}\frac{\partial b_j}{\partial a_k}$$

Given vectors, vector-wise:
$$\boldsymbol{\nabla_a c} = \boldsymbol{\nabla_b c}\, \boldsymbol{\nabla_a b^T}$$

In denominator layout (≠ numerator layout):

$$\boldsymbol{a} = \begin{bmatrix} a_1 \\ \vdots \\ a_{n_a} \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_{n_b} \end{bmatrix} \quad \boldsymbol{c} = \begin{bmatrix} c_1 \\ \vdots \\ c_{n_c} \end{bmatrix} \qquad \boldsymbol{\nabla_b c} = \begin{bmatrix} \dfrac{\partial c_1}{\partial b_1} & \cdots & \dfrac{\partial c_{n_c}}{\partial b_1} \\ \vdots & & \vdots \\ \dfrac{\partial c_1}{\partial b_{n_b}} & \cdots & \dfrac{\partial c_{n_c}}{\partial b_{n_b}} \end{bmatrix}$$

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

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

$$= \sum_j \frac{-1_{y=j}}{\partial \boldsymbol{f}(\boldsymbol{x})_y} \frac{\partial softmax(\tilde{y})_j}{\partial \tilde{y}_i}$$

$$= \begin{cases} \dfrac{-1}{\boldsymbol{f}(\boldsymbol{x})_y} softmax(\tilde{y})_y \big(1 - softmax(\tilde{y})_y\big) & if\ i = y \\[2em] \dfrac{1}{\boldsymbol{f}(\boldsymbol{x})_y} softmax(\tilde{y})_y softmax(\tilde{y})_i & if\ i \neq y \end{cases}$$

$$= \begin{cases} -1 + \boldsymbol{f}(\boldsymbol{x})_y & if\ i = y \\ \boldsymbol{f}(\boldsymbol{x})_i & if\ i \neq y \end{cases}$$

$\color{red}{\nabla_{\tilde{y}}\mathcal{L} = \boldsymbol{f}(\boldsymbol{x}) - \boldsymbol{e}(y)}$ with one-hot encoding of the target



43

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

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

$$= \sum_j \frac{-1_{y=j}}{\partial \boldsymbol{f}(\boldsymbol{x})_y} \frac{\partial softmax(\tilde{y})_j}{\partial \tilde{y}_i}$$
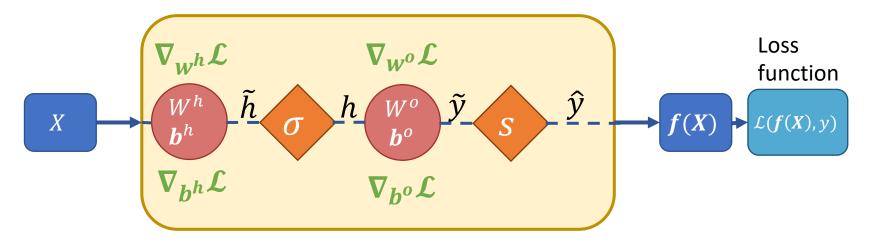
# Backpropagation

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

$$= \sum_j \frac{-1_{y=j}}{\partial \boldsymbol{f(x)}_y} \frac{\partial softmax(\tilde{y})_j}{\partial \tilde{y}_i}$$

$$= \begin{cases} \dfrac{-1}{\boldsymbol{f(x)}_y} softmax(\tilde{y})_y \big(1 - softmax(\tilde{y})_y\big) & if\ i = y \\[2em] \dfrac{1}{\boldsymbol{f(x)}_y} softmax(\tilde{y})_y softmax(\tilde{y})_i & if\ i \neq y \end{cases}$$



46

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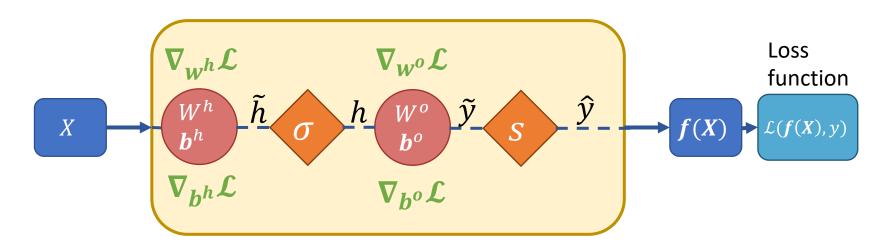
$$\frac{\partial \mathcal{L}}{\partial \tilde{y}_i} = \sum_j \frac{\partial \mathcal{L}}{\partial \boldsymbol{f}(\boldsymbol{x})_j} \frac{\partial \boldsymbol{f}(\boldsymbol{x})_j}{\partial \tilde{y}_i}$$

$$= \sum_j \frac{-1_{y=j}}{\partial \boldsymbol{f}(\boldsymbol{x})_y} \frac{\partial softmax(\tilde{y})_j}{\partial \tilde{y}_i}$$

$$= \begin{cases} \dfrac{-1}{\boldsymbol{f}(\boldsymbol{x})_y} softmax(\tilde{y})_y \big(1 - softmax(\tilde{y})_y\big) & if \ i = y \\[2mm] \dfrac{1}{\boldsymbol{f}(\boldsymbol{x})_y} softmax(\tilde{y})_y softmax(\tilde{y})_i & if \ i \neq y \end{cases}$$

$$= \begin{cases} -1 + \boldsymbol{f}(\boldsymbol{x})_y & if \ i = y \\ \boldsymbol{f}(\boldsymbol{x})_i & if \ i \neq y \end{cases}$$



47

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

$$= \sum_j \frac{-1_{y=j}}{\partial \boldsymbol{f}(\boldsymbol{x})_y} \frac{\partial softmax(\tilde{y})_j}{\partial \tilde{y}_i}$$

$$= \begin{cases} \frac{-1}{\boldsymbol{f}(\boldsymbol{x})_y} softmax(\tilde{y})_y \big(1 - softmax(\tilde{y})_y\big) & if\ i = y \\[2em] \frac{1}{\boldsymbol{f}(\boldsymbol{x})_y} softmax(\tilde{y})_y softmax(\tilde{y})_i & if\ i \neq y \end{cases}$$

$$= \begin{cases} -1 + \boldsymbol{f}(\boldsymbol{x})_y & if\ i = y \\ \boldsymbol{f}(\boldsymbol{x})_i & if\ i \neq y \end{cases}$$

$\textcolor{red}{\nabla_{\tilde{y}} \mathcal{L} = \boldsymbol{f}(\boldsymbol{x}) - \boldsymbol{e}(y) \text{ with one-hot encoding of the target}}$

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

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

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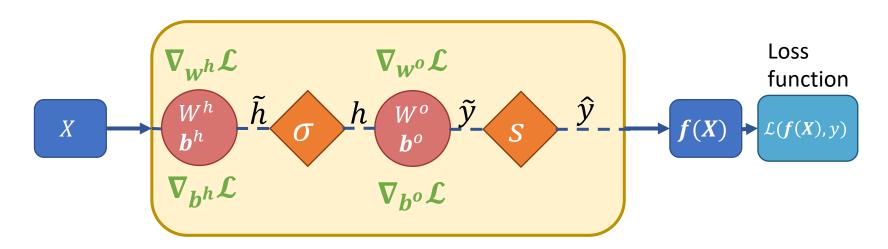
$$\tilde{y} = W^o h + b^o$$

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

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

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

$$\nabla_h\mathcal{L} = W^{o^T} \cdot \nabla_{\tilde{y}}\mathcal{L}$$

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

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

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

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

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

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

$$\nabla_h\mathcal{L} = W^{o^T} \cdot \nabla_{\tilde{y}}\mathcal{L}$$

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

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

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

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

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

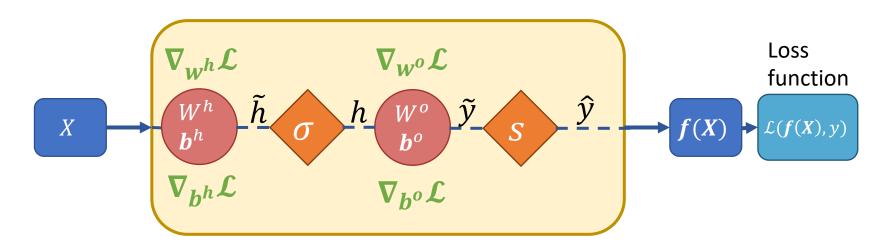$$\nabla_{w^o}\mathcal{L} = \nabla_{\tilde{y}}\mathcal{L} \cdot h^T$$

$$\nabla_h\mathcal{L} = W^{o^T} \cdot \nabla_{\tilde{y}}\mathcal{L}$$

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

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

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

# Backpropagation

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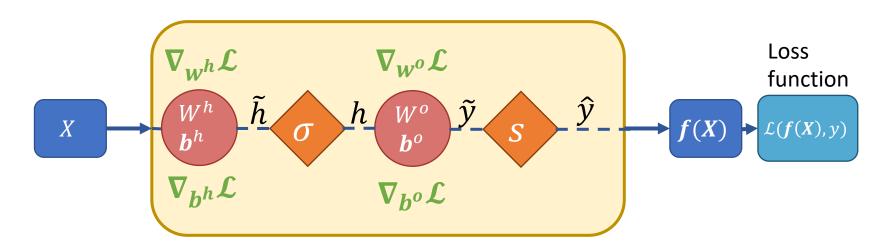
$$\nabla_{b^o}\mathcal{L} = \nabla_{\tilde{y}}\mathcal{L}$$

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

$$\nabla_h\mathcal{L} = W^{o^T} \cdot \nabla_{\tilde{y}}\mathcal{L}$$

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

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

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

When in doubt, look at the shape.

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

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

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

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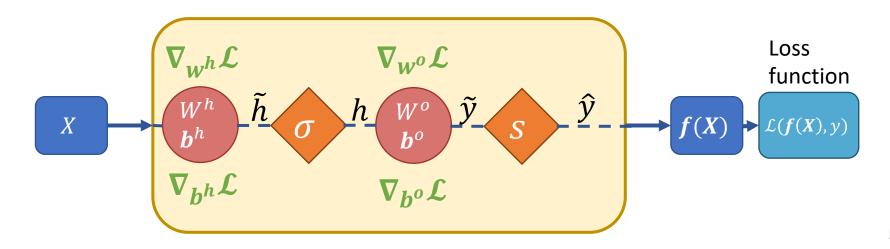
$$\nabla_h\mathcal{L} = W^{o^T} \cdot \nabla_{\tilde{y}}\mathcal{L}$$

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

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

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

To have huge speed-up:
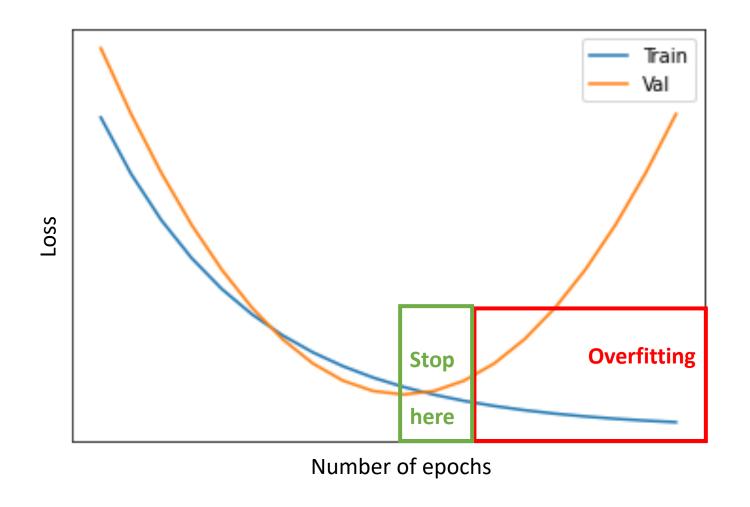1. Re-use previous gradients
2. Save tensors during forward

# Tips & Tricks

# When to Stop?

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

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

# (Mini-)Batch Size

**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->64, lr 0.1->0.2

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

Controls the rate of change.

Too high:
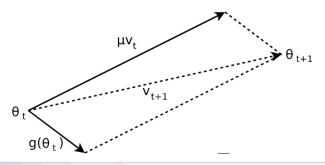→ Cannot converge, but diverge, reduce overfitting

Too low:
→ Super slow, stuck in bad local minima

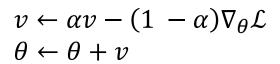Start with high learning rate, and decreases it through time

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

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



Step-size α = 0.02  Momentum β = 0.0

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?
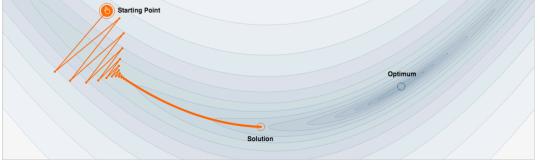


Step-size α = 0.02  Momentum β = 0.85

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?
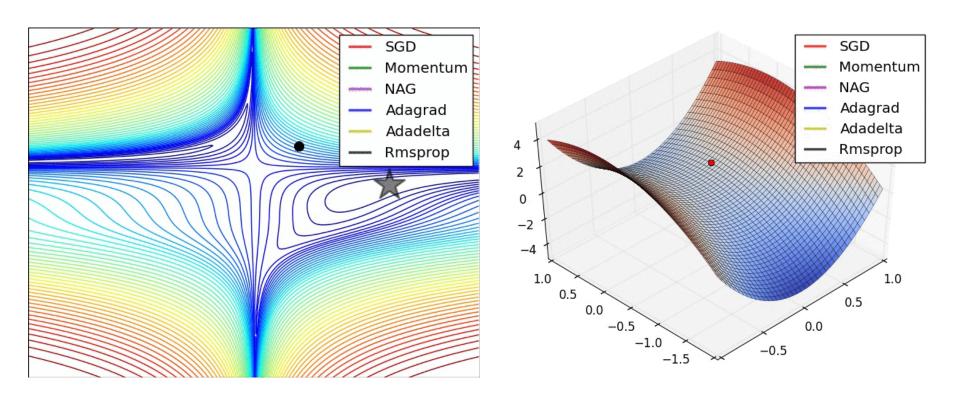
# Optimizers

Modern optimizers have an **adaptive learning rate** per parameter based on gradient statistics.
→ Especially useful on saddle point
→ The most famous is Adam

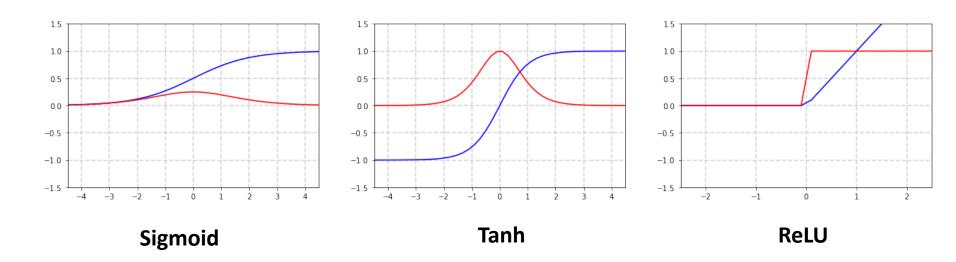But a well-tuned SGD with momentum can sometimes be the best.



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** $LReLU(x) = \begin{cases} x & if\ x > 0 \\ \epsilon\ x & otherwise \end{cases}$

# Initialization

Initializing the biases $b^h$ and $b^o$ to very small values.
$\rightarrow$ Helpful to avoid dying neurons with ReLU

Initializing the weights $W^h$ and $W^o$ to:
- **Zero-weights**
    - $\rightarrow$ No learning because gradient w.r.t input is also zero
- **Constant weights**
    - $\rightarrow$ Symmetry where two hidden neurons are connected to the same inputs, they learn the same pattern!
- **Large values**
    - $\rightarrow$ Risk gradient explosion
- **He / Glorot initialization**
    - $\rightarrow$ Normalize weights to avoid explosion with large number of outgoing connections

# Small break,
# then coding session!