

# Neural networks and Backpropagation

Charles Ollion - Olivier Grisel



# Neural Network for classification

Vector function with tunable parameters  $\theta$

$$\mathbf{f}(\cdot; \theta) : \mathbb{R}^N \rightarrow (0, 1)^K$$

# Neural Network for classification

Vector function with tunable parameters  $\theta$

$$\mathbf{f}(\cdot; \theta) : \mathbb{R}^N \rightarrow (0, 1)^K$$

Sample  $s$  in dataset  $S$ :

- input:  $\mathbf{x}^s \in \mathbb{R}^N$
- expected output:  $y^s \in [0, K - 1]$

# Neural Network for classification

Vector function with tunable parameters  $\theta$

$$\mathbf{f}(\cdot; \theta) : \mathbb{R}^N \rightarrow (0, 1)^K$$

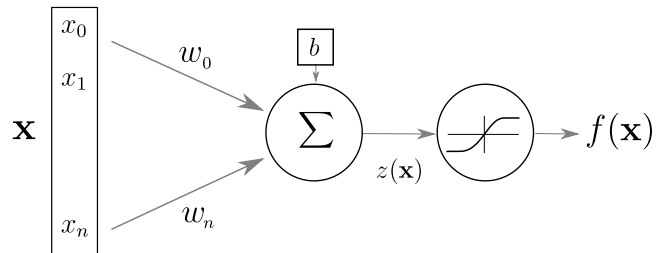
Sample  $s$  in dataset  $S$ :

- input:  $\mathbf{x}^s \in \mathbb{R}^N$
- expected output:  $y^s \in [0, K - 1]$

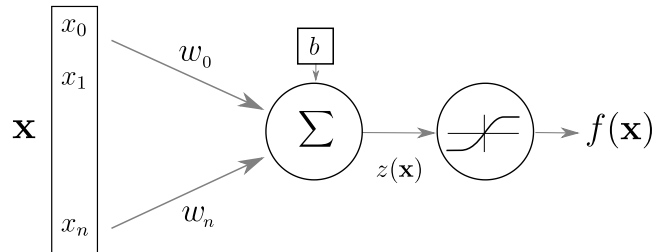
Output is a conditional probability distribution:

$$\mathbf{f}(\mathbf{x}^s; \theta)_c = P(Y = c | X = \mathbf{x}^s)$$

# Artificial Neuron



# Artificial Neuron

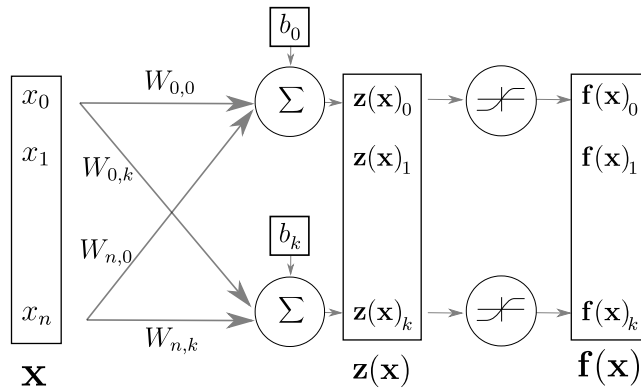


$$z(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

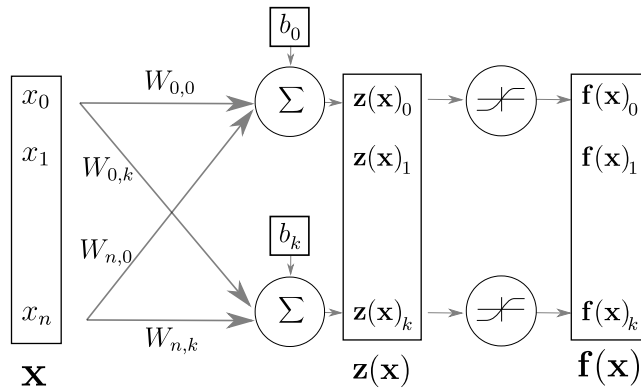
$$f(\mathbf{x}) = g(\mathbf{w}^T \mathbf{x} + b)$$

- $\mathbf{x}, f(\mathbf{x})$  input and output
- $z(\mathbf{x})$  pre-activation
- $\mathbf{w}, b$  weights and bias
- $g$  activation function

# Layer of Neurons



# Layer of Neurons

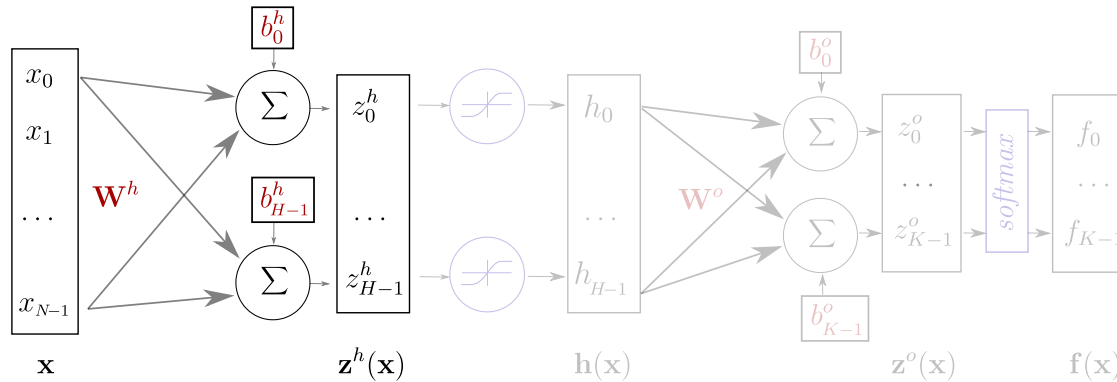


$$\mathbf{f}(\mathbf{x}) = g(\mathbf{z}(\mathbf{x})) = g(\mathbf{W}\mathbf{x} + \mathbf{b})$$

- $\mathbf{W}, \mathbf{b}$  now matrix and vector

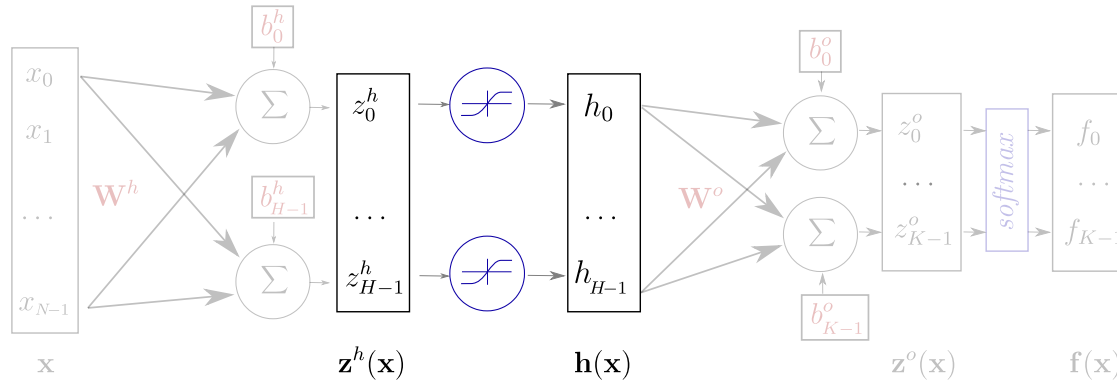


# One Hidden Layer Network



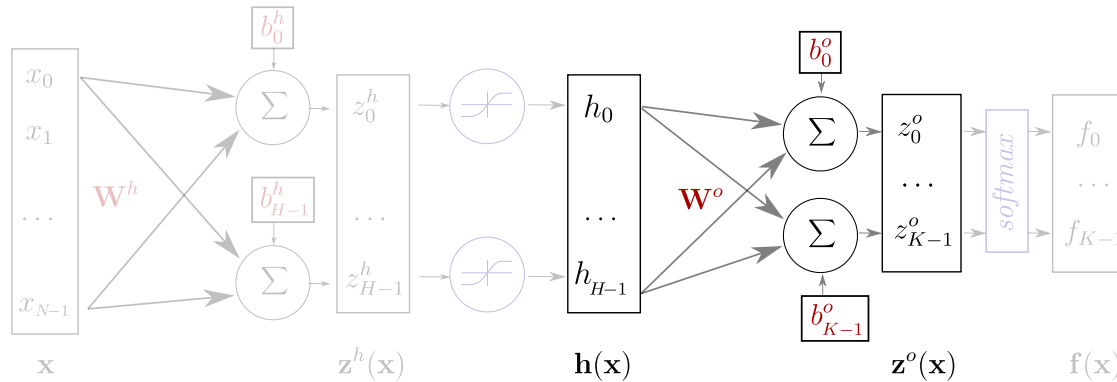
- $\mathbf{z}^h(\mathbf{x}) = \mathbf{W}^h \mathbf{x} + \mathbf{b}^h$
- $\mathbf{h}(\mathbf{x}) = g(\mathbf{z}^h(\mathbf{x})) = g(\mathbf{W}^h \mathbf{x} + \mathbf{b}^h)$
- $\mathbf{z}^o(\mathbf{x}) = \mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o$
- $\mathbf{f}(\mathbf{x}) = \text{softmax}(\mathbf{z}^o) = \text{softmax}(\mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o)$

# One Hidden Layer Network



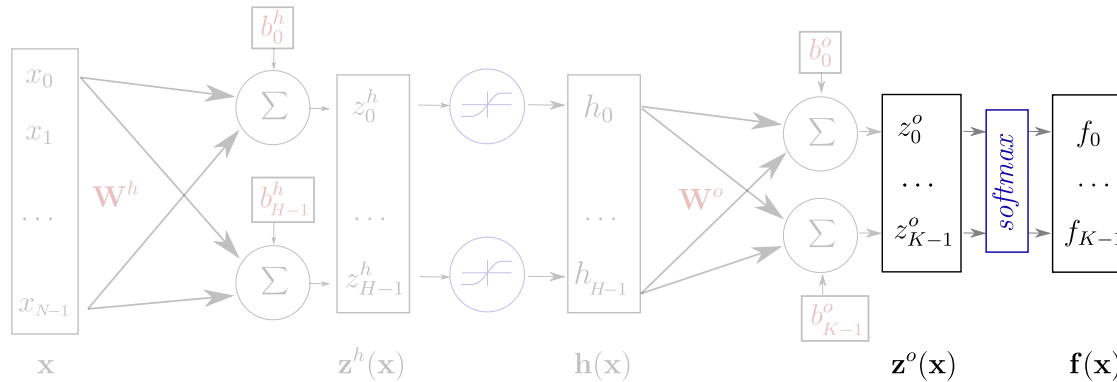
- $\mathbf{z}^h(\mathbf{x}) = \mathbf{W}^h \mathbf{x} + \mathbf{b}^h$
- $\mathbf{h}(\mathbf{x}) = g(\mathbf{z}^h(\mathbf{x})) = g(\mathbf{W}^h \mathbf{x} + \mathbf{b}^h)$
- $\mathbf{z}^o(\mathbf{x}) = \mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o$
- $\mathbf{f}(\mathbf{x}) = \text{softmax}(\mathbf{z}^o) = \text{softmax}(\mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o)$

# One Hidden Layer Network



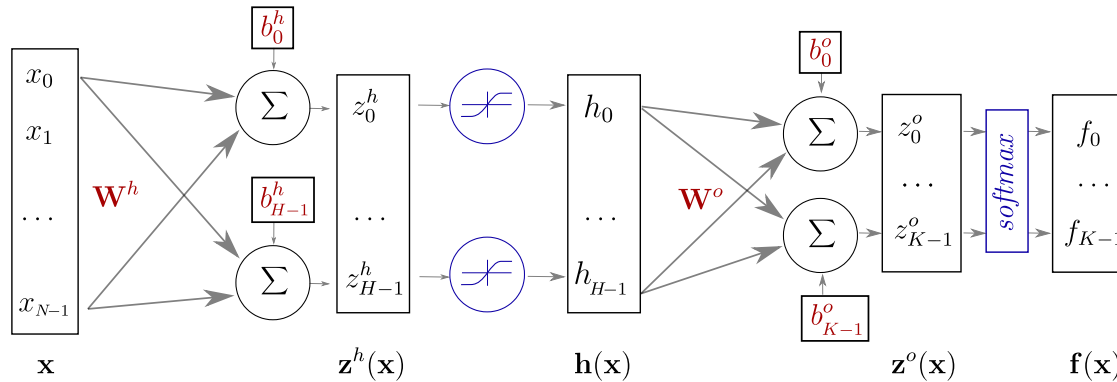
- $\mathbf{z}^h(\mathbf{x}) = \mathbf{W}^h \mathbf{x} + \mathbf{b}^h$
- $\mathbf{h}(\mathbf{x}) = g(\mathbf{z}^h(\mathbf{x})) = g(\mathbf{W}^h \mathbf{x} + \mathbf{b}^h)$
- $\mathbf{z}^o(\mathbf{x}) = \mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o$
- $\mathbf{f}(\mathbf{x}) = \text{softmax}(\mathbf{z}^o) = \text{softmax}(\mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o)$

# One Hidden Layer Network

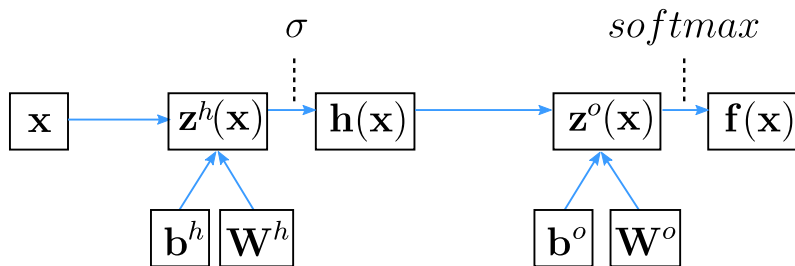


- $\mathbf{z}^h(\mathbf{x}) = \mathbf{W}^h \mathbf{x} + \mathbf{b}^h$
- $\mathbf{h}(\mathbf{x}) = g(\mathbf{z}^h(\mathbf{x})) = g(\mathbf{W}^h \mathbf{x} + \mathbf{b}^h)$
- $\mathbf{z}^o(\mathbf{x}) = \mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o$
- $\mathbf{f}(\mathbf{x}) = \text{softmax}(\mathbf{z}^o) = \text{softmax}(\mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o)$

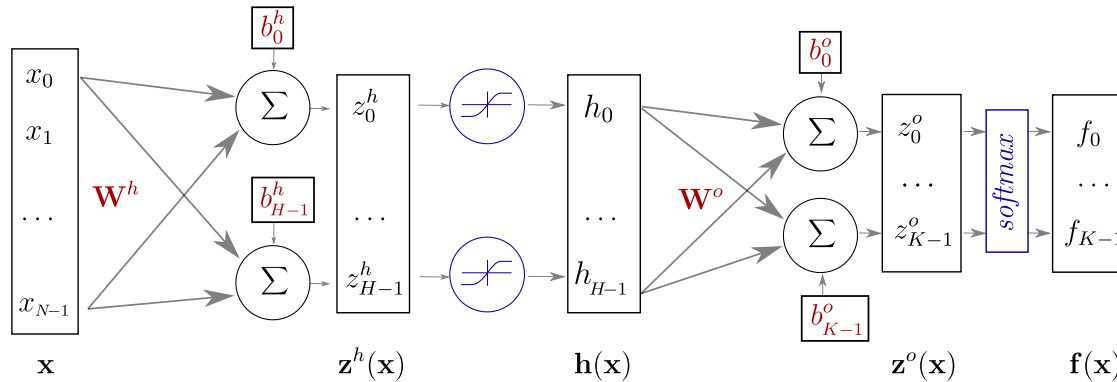
# One Hidden Layer Network



## Alternate representation



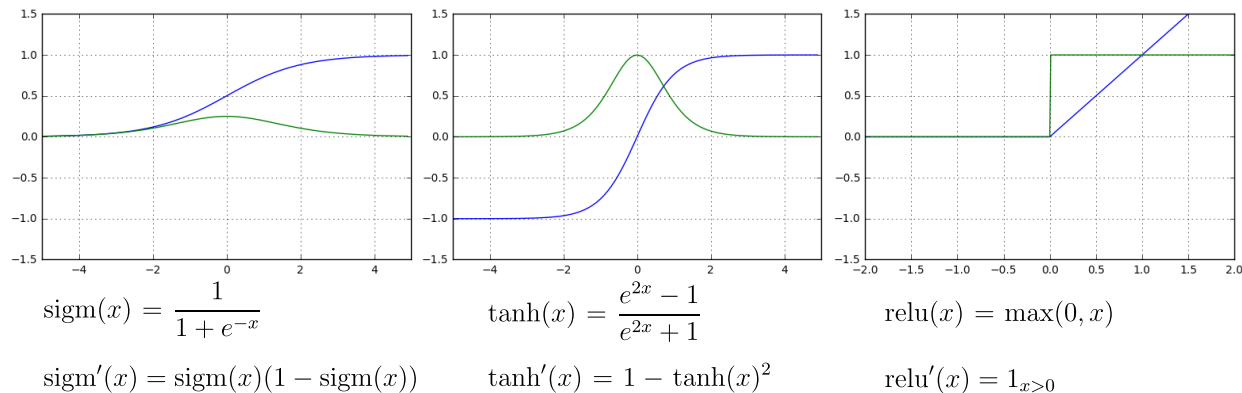
# One Hidden Layer Network



## Keras implementation

```
model = Sequential()
model.add(Dense(H, input_dim=N)) # weight matrix dim [N * H]
model.add(Activation("tanh"))
model.add(Dense(K)) # weight matrix dim [H x K]
model.add(Activation("softmax"))
```

# Element-wise activation functions



- blue: activation function
- green: derivative

# Softmax function

$$\text{softmax}(\mathbf{x}) = \frac{1}{\sum_{i=1}^n e^{x_i}} \cdot \begin{bmatrix} e^{x_1} \\ e^{x_2} \\ \vdots \\ e^{x_n} \end{bmatrix}$$

$$\frac{\partial \text{softmax}(\mathbf{x})_i}{\partial x_j} = \begin{cases} \text{softmax}(\mathbf{x})_i \cdot (1 - \text{softmax}(\mathbf{x})_i) & i = j \\ -\text{softmax}(\mathbf{x})_i \cdot \text{softmax}(\mathbf{x})_j & i \neq j \end{cases}$$



# Softmax function

$$\text{softmax}(\mathbf{x}) = \frac{1}{\sum_{i=1}^n e^{x_i}} \cdot \begin{bmatrix} e^{x_1} \\ e^{x_2} \\ \vdots \\ e^{x_n} \end{bmatrix}$$

$$\frac{\partial \text{softmax}(\mathbf{x})_i}{\partial x_j} = \begin{cases} \text{softmax}(\mathbf{x})_i \cdot (1 - \text{softmax}(\mathbf{x})_i) & i = j \\ -\text{softmax}(\mathbf{x})_i \cdot \text{softmax}(\mathbf{x})_j & i \neq j \end{cases}$$

- vector of values in (0, 1) that add up to 1
- $p(Y = c|X = \mathbf{x}) = \text{softmax}(\mathbf{z}(\mathbf{x}))_c$
- the pre-activation vector  $\mathbf{z}(\mathbf{x})$  is often called "the logits"

# Training the network

Find parameters  $\theta = (\mathbf{W}^h; \mathbf{b}^h; \mathbf{W}^o; \mathbf{b}^o)$  that minimize the **negative log likelihood** (or [cross entropy](#))

# Training the network

Find parameters  $\theta = (\mathbf{W}^h; \mathbf{b}^h; \mathbf{W}^o; \mathbf{b}^o)$  that minimize the **negative log likelihood** (or [cross entropy](#))

The loss function for a given sample  $s \in S$ :

$$l(\mathbf{f}(\mathbf{x}^s; \theta), y^s) = nll(\mathbf{x}^s, y^s; \theta) = -\log \mathbf{f}(\mathbf{x}^s; \theta)_{y^s}$$

# Training the network

Find parameters  $\theta = (\mathbf{W}^h; \mathbf{b}^h; \mathbf{W}^o; \mathbf{b}^o)$  that minimize the **negative log likelihood** (or [cross entropy](#))

The loss function for a given sample  $s \in S$ :

$$l(\mathbf{f}(\mathbf{x}^s; \theta), y^s) = nll(\mathbf{x}^s, y^s; \theta) = -\log \mathbf{f}(\mathbf{x}^s; \theta)_{y^s}$$

example  $y^s = 3$

$$l(\mathbf{f}(\mathbf{x}^s; \theta), y^s) = l \left( \begin{bmatrix} f_0 \\ \dots \\ f_3 \\ \dots \\ f_{K-1} \end{bmatrix}, \begin{bmatrix} 0 \\ \dots \\ 1 \\ \dots \\ 0 \end{bmatrix} \right) = -\log f_3$$

# Training the network

Find parameters  $\theta = (\mathbf{W}^h; \mathbf{b}^h; \mathbf{W}^o; \mathbf{b}^o)$  that minimize the **negative log likelihood** (or [cross entropy](#))

The loss function for a given sample  $s \in S$ :

$$l(\mathbf{f}(\mathbf{x}^s; \theta), y^s) = nll(\mathbf{x}^s, y^s; \theta) = -\log \mathbf{f}(\mathbf{x}^s; \theta)_{y^s}$$

The cost function is the negative likelihood of the model computed on the full training set (for i.i.d. samples):

$$L_S(\theta) = -\frac{1}{|S|} \sum_{s \in S} \log \mathbf{f}(\mathbf{x}^s; \theta)_{y^s}$$

# Training the network

Find parameters  $\theta = (\mathbf{W}^h; \mathbf{b}^h; \mathbf{W}^o; \mathbf{b}^o)$  that minimize the **negative log likelihood** (or [cross entropy](#))

The loss function for a given sample  $s \in S$ :

$$l(\mathbf{f}(\mathbf{x}^s; \theta), y^s) = nll(\mathbf{x}^s, y^s; \theta) = -\log \mathbf{f}(\mathbf{x}^s; \theta)_{y^s}$$

The cost function is the negative likelihood of the model computed on the full training set (for i.i.d. samples):

$$L_S(\theta) = -\frac{1}{|S|} \sum_{s \in S} \log \mathbf{f}(\mathbf{x}^s; \theta)_{y^s} + \lambda \Omega(\theta)$$

$\lambda \Omega(\theta) = \lambda(\|\mathbf{W}^h\|^2 + \|\mathbf{W}^o\|^2)$  is an optional regularization term.

# Stochastic Gradient Descent

Initialize  $\theta$  randomly

# Stochastic Gradient Descent

Initialize  $\theta$  randomly

For  $E$  epochs perform:

- Randomly select a small batch of samples ( $B \subset S$ )



# Stochastic Gradient Descent

Initialize  $\theta$  randomly

For  $E$  epochs perform:

- Randomly select a small batch of samples ( $B \subset S$ )
  - Compute gradients:  $\Delta = \nabla_{\theta} L_B(\theta)$

# Stochastic Gradient Descent

Initialize  $\theta$  randomly

For  $E$  epochs perform:

- Randomly select a small batch of samples ( $B \subset S$ )
  - Compute gradients:  $\Delta = \nabla_{\theta} L_B(\theta)$
  - Update parameters:  $\theta \leftarrow \theta - \eta \Delta$
  - $\eta > 0$  is called the learning rate

# Stochastic Gradient Descent

Initialize  $\theta$  randomly

For  $E$  epochs perform:

- Randomly select a small batch of samples ( $B \subset S$ )
  - Compute gradients:  $\Delta = \nabla_{\theta} L_B(\theta)$
  - Update parameters:  $\theta \leftarrow \theta - \eta \Delta$
  - $\eta > 0$  is called the learning rate
- Repeat until the epoch is completed (all of  $S$  is covered)

# Stochastic Gradient Descent

Initialize  $\theta$  randomly

For  $E$  epochs perform:

- Randomly select a small batch of samples ( $B \subset S$ )
  - Compute gradients:  $\Delta = \nabla_{\theta} L_B(\theta)$
  - Update parameters:  $\theta \leftarrow \theta - \eta \Delta$
  - $\eta > 0$  is called the learning rate
- Repeat until the epoch is completed (all of  $S$  is covered)

Stop when reaching criterion:

- nll stops decreasing when computed on validation set

# Computing Gradients

Output Weights:  $\frac{\partial l(\mathbf{f}(\mathbf{x}), y)}{\partial W_{i,j}^o}$

Output bias:  $\frac{\partial l(\mathbf{f}(\mathbf{x}), y)}{\partial b_i^o}$

Hidden Weights:  $\frac{\partial l(\mathbf{f}(\mathbf{x}), y)}{\partial W_{i,j}^h}$

Hidden bias:  $\frac{\partial l(\mathbf{f}(\mathbf{x}), y)}{\partial b_i^h}$

# Computing Gradients

Output Weights:  $\frac{\partial l(\mathbf{f}(\mathbf{x}), y)}{\partial W_{i,j}^o}$

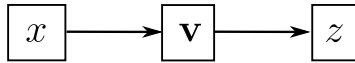
Output bias:  $\frac{\partial l(\mathbf{f}(\mathbf{x}), y)}{\partial b_i^o}$

Hidden Weights:  $\frac{\partial l(\mathbf{f}(\mathbf{x}), y)}{\partial W_{i,j}^h}$

Hidden bias:  $\frac{\partial l(\mathbf{f}(\mathbf{x}), y)}{\partial b_i^h}$

- The network is a composition of differentiable modules
- We can apply the "chain rule"

# Chain rule



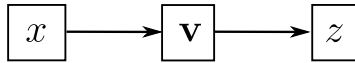
$$z = u(\mathbf{v}(x))$$

$$\frac{\partial z}{\partial x} = ?$$

$$v_j$$

$\mathbf{v}$

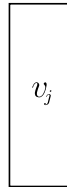
# Chain rule



$$z = u(\mathbf{v}(x))$$

**chain-rule**

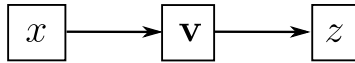
$$\frac{\partial z}{\partial x} = \sum_j \frac{\partial z}{\partial v_j} \frac{\partial v_j}{\partial x}$$



**v**



# Chain rule



$$z = u(\mathbf{v}(x))$$

**chain-rule**

$$\frac{\partial z}{\partial x} = \sum_j \frac{\partial z}{\partial v_j} \frac{\partial v_j}{\partial x} = \nabla u \cdot \frac{\partial \mathbf{v}}{\partial x}$$

$$v_j$$

$\mathbf{v}$

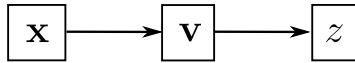
$$\frac{\partial v_j}{\partial x}$$

$\frac{\partial \mathbf{v}}{\partial x}$

$$\frac{\partial z}{\partial v_j}$$

$\nabla u$

# Chain rule



$$z = u(\mathbf{v}(\mathbf{x}))$$

**chain-rule**

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial v_j} \frac{\partial v_j}{\partial x_i} = \nabla u \cdot \frac{\partial \mathbf{v}}{\partial x_i}$$

$$v_j$$

$\mathbf{v}$

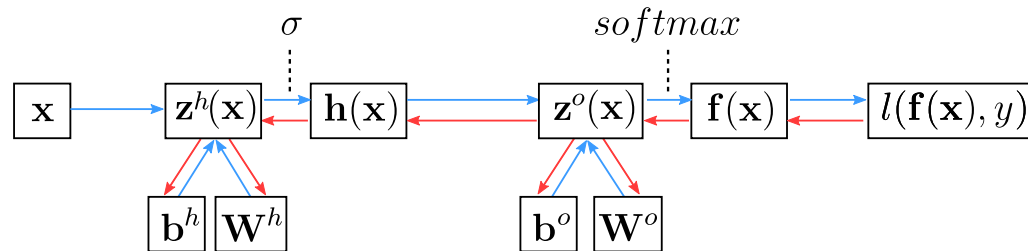
$$\frac{\partial v_j}{\partial x_i}$$

$$\frac{\partial \mathbf{v}}{\partial x_i}$$

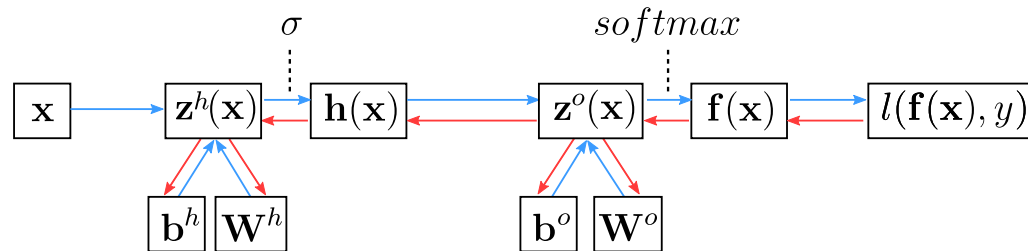
$$\frac{\partial z}{\partial v_j}$$

$$\nabla u$$

# Backpropagation



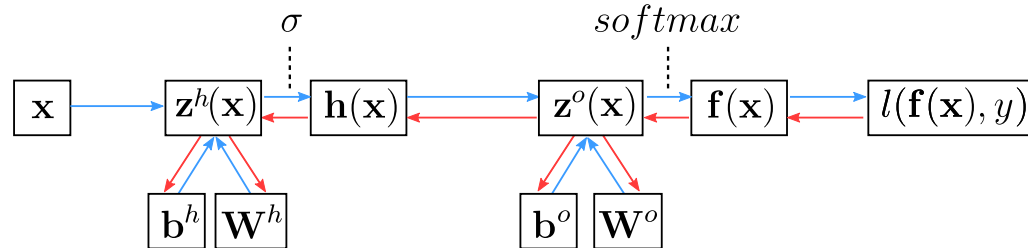
# Backpropagation



Compute partial derivatives of the loss

- $$\frac{\partial l(\mathbf{f}(\mathbf{x}), y)}{\partial \mathbf{f}(\mathbf{x})_i} = \frac{\partial -\log \mathbf{f}(\mathbf{x})_y}{\partial \mathbf{f}(\mathbf{x})_i} = \frac{-1_{y=i}}{\mathbf{f}(\mathbf{x})_y} = \frac{\partial l}{\partial \mathbf{f}(\mathbf{x})_i}$$

# Backpropagation

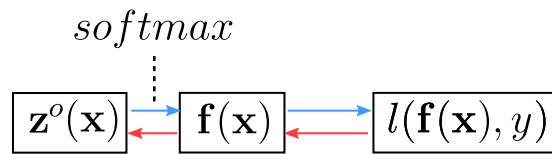


Compute partial derivatives of the loss

- $\frac{\partial l(\mathbf{f}(\mathbf{x}), y)}{\partial \mathbf{f}(\mathbf{x})_i} = \frac{\partial -\log \mathbf{f}(\mathbf{x})_y}{\partial \mathbf{f}(\mathbf{x})_i} = \frac{-1_{y=i}}{\mathbf{f}(\mathbf{x})_y} = \frac{\partial l}{\partial \mathbf{f}(\mathbf{x})_i}$
- $\frac{\partial l}{\partial \mathbf{z}^o(\mathbf{x})_i} = ?$

$$\frac{\partial l}{\partial \mathbf{z}^o(\mathbf{x})_i} = \sum_j \frac{\partial l}{\partial \mathbf{f}(\mathbf{x})_j} \frac{\partial \mathbf{f}(\mathbf{x})_j}{\partial \mathbf{z}^o(\mathbf{x})_i}$$

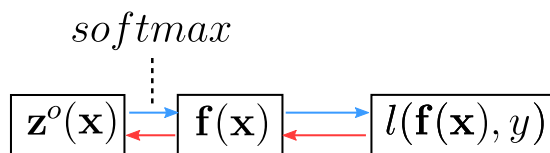
Chain rule!



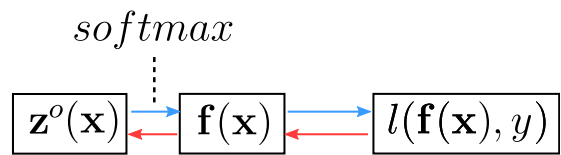
$$\begin{aligned}\frac{\partial l}{\partial \mathbf{z}^o(\mathbf{x})_i} &= \sum_j \frac{\partial l}{\partial \mathbf{f}(\mathbf{x})_j} \frac{\partial \mathbf{f}(\mathbf{x})_j}{\partial \mathbf{z}^o(\mathbf{x})_i} \\ &= \sum_j \frac{-1_{y=j}}{\mathbf{f}(\mathbf{x})_y} \frac{\partial \text{softmax}(\mathbf{z}^o(\mathbf{x}))_j}{\partial \mathbf{z}^o(\mathbf{x})_i}\end{aligned}$$

$$\frac{\partial l}{\partial \mathbf{f}(\mathbf{x})_i} = \frac{-1_{y=i}}{\mathbf{f}(\mathbf{x})_y}$$

$$\mathbf{f}(\mathbf{x}) = \text{softmax}(\mathbf{z}^o(\mathbf{x}))$$

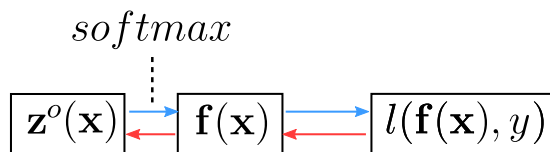


$$\begin{aligned}
\frac{\partial l}{\partial \mathbf{z}^o(\mathbf{x})_i} &= \sum_j \frac{\partial l}{\partial \mathbf{f}(\mathbf{x})_j} \frac{\partial \mathbf{f}(\mathbf{x})_j}{\partial \mathbf{z}^o(\mathbf{x})_i} \\
&= \sum_j \frac{-1_{y=j}}{\mathbf{f}(\mathbf{x})_y} \frac{\partial \text{softmax}(\mathbf{z}^o(\mathbf{x}))_j}{\partial \mathbf{z}^o(\mathbf{x})_i} \\
&= -\frac{1}{\mathbf{f}(\mathbf{x})_y} \frac{\partial \text{softmax}(\mathbf{z}^o(\mathbf{x}))_y}{\partial \mathbf{z}^o(\mathbf{x})_i}
\end{aligned}$$





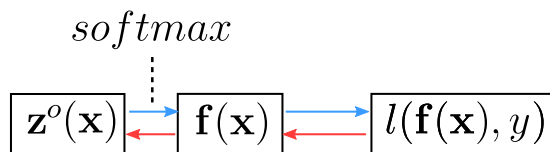
$$\begin{aligned}
\frac{\partial l}{\partial \mathbf{z}^o(\mathbf{x})_i} &= \sum_j \frac{\partial l}{\partial \mathbf{f}(\mathbf{x})_j} \frac{\partial \mathbf{f}(\mathbf{x})_j}{\partial \mathbf{z}^o(\mathbf{x})_i} \\
&= \sum_j \frac{-1_{y=j}}{\mathbf{f}(\mathbf{x})_y} \frac{\partial \text{softmax}(\mathbf{z}^o(\mathbf{x}))_j}{\partial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



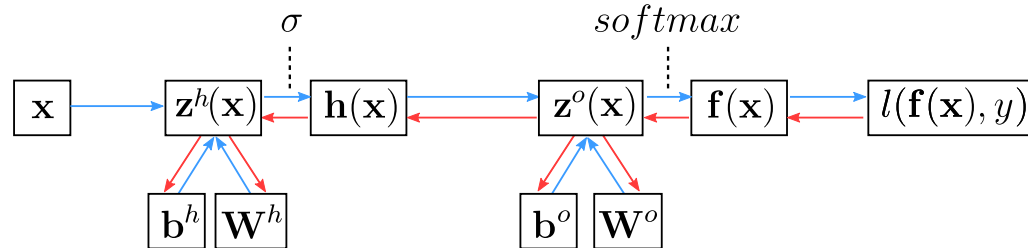
$$\begin{aligned}
\frac{\partial l}{\partial \mathbf{z}^o(\mathbf{x})_i} &= \sum_j \frac{\partial l}{\partial \mathbf{f}(\mathbf{x})_j} \frac{\partial \mathbf{f}(\mathbf{x})_j}{\partial \mathbf{z}^o(\mathbf{x})_i} \\
&= \sum_j \frac{-1_{y=j}}{\mathbf{f}(\mathbf{x})_y} \frac{\partial \text{softmax}(\mathbf{z}^o(\mathbf{x}))_j}{\partial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

$$\nabla_{\mathbf{z}^o(\mathbf{x})} l(\mathbf{f}(\mathbf{x}), y) = \mathbf{f}(\mathbf{x}) - \mathbf{e}(y)$$

$\mathbf{e}(y)$ : one-hot encoding of  $y$



# Backpropagation

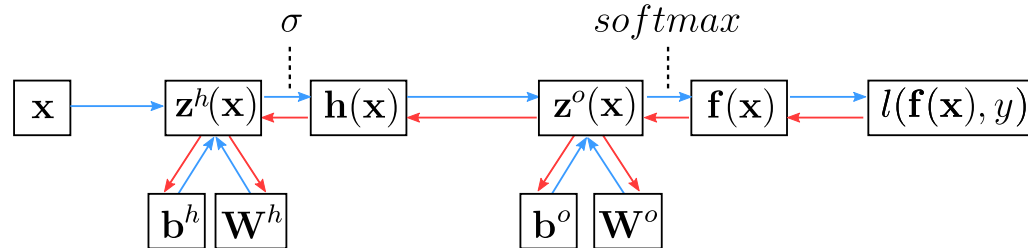


## Gradients

- $\nabla_{\mathbf{z}^o(\mathbf{x})} l = \mathbf{f}(\mathbf{x}) - \mathbf{e}(y)$
- $\nabla_{\mathbf{b}^o} l = \mathbf{f}(\mathbf{x}) - \mathbf{e}(y)$

because  $\mathbf{z}^o(\mathbf{x}) = \mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o$  and then  $\frac{\partial \mathbf{z}^o(\mathbf{x})_i}{\partial \mathbf{b}_j^o} = 1_{i=j}$

# Backpropagation



Partial derivatives related to  $\mathbf{W}^o$

- $\frac{\partial l}{\partial W_{i,j}^o} = \sum_k \frac{\partial l}{\partial \mathbf{z}^o(\mathbf{x})_k} \frac{\partial \mathbf{z}^o(\mathbf{x})_k}{\partial W_{i,j}^o}$
- $\nabla_{\mathbf{W}^o} l = (\mathbf{f}(\mathbf{x}) - \mathbf{e}(y)) \cdot \mathbf{h}(\mathbf{x})^\top$

# Backprop gradients

Compute activation gradients

- $\nabla_{\mathbf{z}^o(\mathbf{x})} l = \mathbf{f}(\mathbf{x}) - \mathbf{e}(y)$

# Backprop gradients

## Compute activation gradients

- $\nabla_{\mathbf{z}^o(\mathbf{x})} \mathbf{l} = \mathbf{f}(\mathbf{x}) - \mathbf{e}(y)$

## Compute layer params gradients

- $\nabla_{\mathbf{W}^o} \mathbf{l} = \nabla_{\mathbf{z}^o(\mathbf{x})} \mathbf{l} \cdot \mathbf{h}(\mathbf{x})^\top$
- $\nabla_{\mathbf{b}^o} \mathbf{l} = \nabla_{\mathbf{z}^o(\mathbf{x})} \mathbf{l}$

# Backprop gradients

## Compute activation gradients

- $\nabla_{\mathbf{z}^o(\mathbf{x})} \mathbf{l} = \mathbf{f}(\mathbf{x}) - \mathbf{e}(y)$

## Compute layer params gradients

- $\nabla_{\mathbf{W}^o} \mathbf{l} = \nabla_{\mathbf{z}^o(\mathbf{x})} \mathbf{l} \cdot \mathbf{h}(\mathbf{x})^\top$
- $\nabla_{\mathbf{b}^o} \mathbf{l} = \nabla_{\mathbf{z}^o(\mathbf{x})} \mathbf{l}$

## Compute prev layer activation gradients

- $\nabla_{\mathbf{h}(\mathbf{x})} \mathbf{l} = \mathbf{W}^{o\top} \nabla_{\mathbf{z}^o(\mathbf{x})} \mathbf{l}$
- $\nabla_{\mathbf{z}^h(\mathbf{x})} \mathbf{l} = \nabla_{\mathbf{h}(\mathbf{x})} \mathbf{l} \odot \sigma'(\mathbf{z}^h(\mathbf{x}))$

# Loss, Initialization and Learning Tricks



# Discrete output (classification)

- Binary classification:  $y \in [0, 1]$ 
  - $Y|X = \mathbf{x} \sim \text{Bernoulli}(b = f(\mathbf{x}; \theta))$
  - output function:  $\text{logistic}(x) = \frac{1}{1+e^{-x}}$
  - loss function: binary cross-entropy
- Multiclass classification:  $y \in [0, K - 1]$ 
  - $Y|X = \mathbf{x} \sim \text{Multinoulli}(\mathbf{p} = \mathbf{f}(\mathbf{x}; \theta))$
  - output function: *softmax*
  - loss function: categorical cross-entropy

# Continuous output (regression)

- Continuous output:  $\mathbf{y} \in \mathbb{R}^n$ 
  - $Y|X = \mathbf{x} \sim \mathcal{N}(\mu = \mathbf{f}(\mathbf{x}; \theta), \sigma^2 \mathbf{I})$
  - output function: Identity
  - loss function: square loss
- Heteroschedastic if  $\mathbf{f}(\mathbf{x}; \theta)$  predicts both  $\mu$  and  $\sigma^2$
- Mixture Density Network (multimodal output)
  - $Y|X = \mathbf{x} \sim GMM_{\mathbf{x}}$
  - $\mathbf{f}(\mathbf{x}; \theta)$  predicts all the parameters: the means, covariance matrices and mixture weights

# Initialization and normalization

- Input data should be normalized to have approx. same range:
  - standardization or quantile normalization

# Initialization and normalization

- Input data should be normalized to have approx. same range:
  - standardization or quantile normalization
- Initializing  $W^h$  and  $W^o$ :
  - Zero is a saddle point: no gradient, no learning

# Initialization and normalization

- Input data should be normalized to have approx. same range:
  - standardization or quantile normalization
- Initializing  $W^h$  and  $W^o$ :
  - Zero is a saddle point: no gradient, no learning
  - Constant init: hidden units collapse by symmetry

# Initialization and normalization

- Input data should be normalized to have approx. same range:
  - standardization or quantile normalization
- Initializing  $W^h$  and  $W^o$ :
  - Zero is a saddle point: no gradient, no learning
  - Constant init: hidden units collapse by symmetry
  - Solution: random init, ex:  $w \sim \mathcal{N}(0, 0.01)$

# Initialization and normalization

- Input data should be normalized to have approx. same range:
  - standardization or quantile normalization
- Initializing  $W^h$  and  $W^o$ :
  - Zero is a saddle point: no gradient, no learning
  - Constant init: hidden units collapse by symmetry
  - Solution: random init, ex:  $w \sim \mathcal{N}(0, 0.01)$
  - Better inits: Xavier Glorot and Kaming He & orthogonal

# Initialization and normalization

- Input data should be normalized to have approx. same range:
  - standardization or quantile normalization
- Initializing  $W^h$  and  $W^o$ :
  - Zero is a saddle point: no gradient, no learning
  - Constant init: hidden units collapse by symmetry
  - Solution: random init, ex:  $w \sim \mathcal{N}(0, 0.01)$
  - Better inits: Xavier Glorot and Kaming He & orthogonal
- Biases can (should) be initialized to zero



# SGD learning rate

- Very sensitive:
  - Too high  $\rightarrow$  early plateau or even divergence
  - Too low  $\rightarrow$  slow convergence

# SGD learning rate

- Very sensitive:
  - Too high  $\rightarrow$  early plateau or even divergence
  - Too low  $\rightarrow$  slow convergence
  - Try a large value first:  $\eta = 0.1$  or even  $\eta = 1$
  - Divide by 10 and retry in case of divergence

# SGD learning rate

- Very sensitive:
  - Too high  $\rightarrow$  early plateau or even divergence
  - Too low  $\rightarrow$  slow convergence
  - Try a large value first:  $\eta = 0.1$  or even  $\eta = 1$
  - Divide by 10 and retry in case of divergence
- Large constant LR prevents final convergence
  - multiply  $\eta_t$  by  $\beta < 1$  after each update

# SGD learning rate

- Very sensitive:
  - Too high  $\rightarrow$  early plateau or even divergence
  - Too low  $\rightarrow$  slow convergence
  - Try a large value first:  $\eta = 0.1$  or even  $\eta = 1$
  - Divide by 10 and retry in case of divergence
- Large constant LR prevents final convergence
  - multiply  $\eta_t$  by  $\beta < 1$  after each update
  - or monitor validation loss and divide  $\eta_t$  by 2 or 10 when no progress
  - See [ReduceLROnPlateau](#) in Keras

# Momentum

Accumulate gradients across successive updates:

$$m_t = \gamma m_{t-1} + \eta \nabla_{\theta} L_{B_t}(\theta_{t-1})$$

$$\theta_t = \theta_{t-1} - m_t$$

$\gamma$  is typically set to 0.9

# Momentum

Accumulate gradients across successive updates:

$$m_t = \gamma m_{t-1} + \eta \nabla_{\theta} L_{B_t}(\theta_{t-1})$$
$$\theta_t = \theta_{t-1} - m_t$$

$\gamma$  is typically set to 0.9

Larger updates in directions where the gradient sign is constant to accelerate in low curvature areas

# Momentum

Accumulate gradients across successive updates:

$$m_t = \gamma m_{t-1} + \eta \nabla_{\theta} L_{B_t}(\theta_{t-1})$$
$$\theta_t = \theta_{t-1} - m_t$$

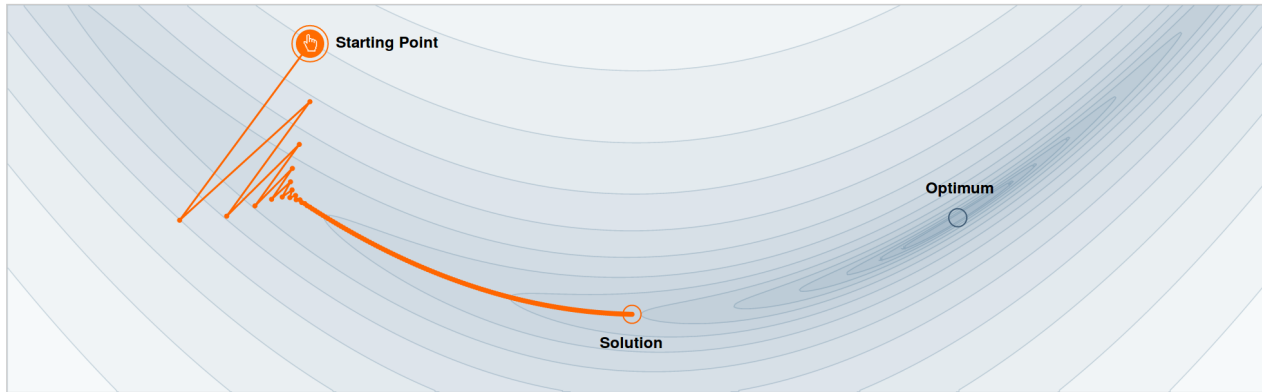
$\gamma$  is typically set to 0.9

Larger updates in directions where the gradient sign is constant to accelerate in low curvature areas

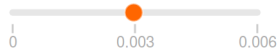
## Nesterov accelerated gradient

$$m_t = \gamma m_{t-1} + \eta \nabla_{\theta} L_{B_t}(\theta_{t-1} - \gamma m_{t-1})$$
$$\theta_t = \theta_{t-1} - m_t$$

Better at handling changes in gradient direction.



Step-size  $\alpha = 0.0030$



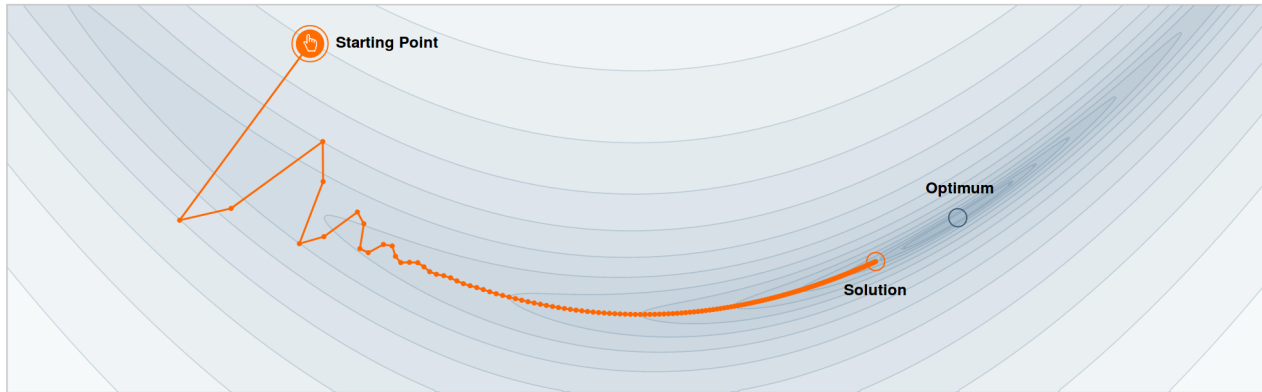
Momentum  $\beta = 0.0$



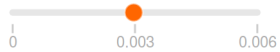
We often think of Momentum as a means of dampening oscillations and speeding up the iterations, leading to faster convergence. But it has other interesting behavior. It allows a larger range of step-sizes to be used, and creates its own oscillations. What is going on?

## [Why Momentum Really Works](#)





Step-size  $\alpha = 0.0030$

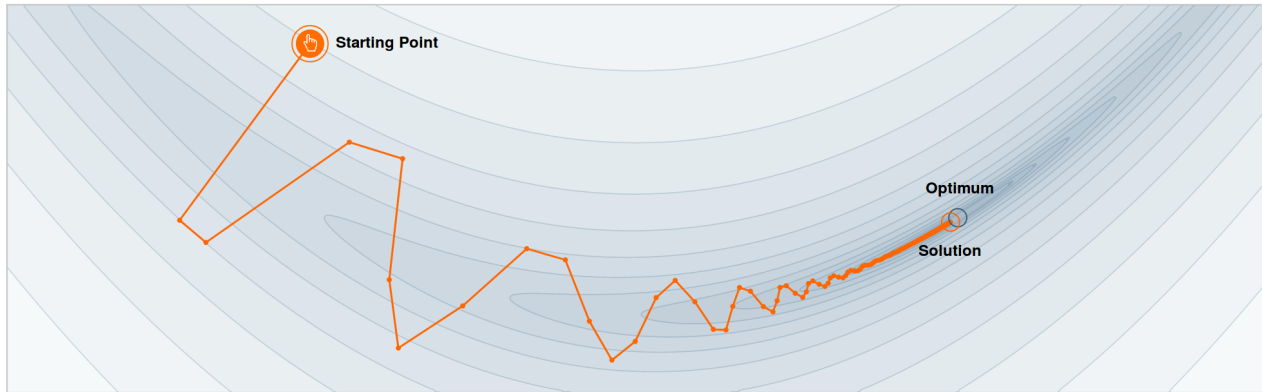


Momentum  $\beta = 0.60$

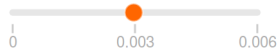


We often think of Momentum as a means of dampening oscillations and speeding up the iterations, leading to faster convergence. But it has other interesting behavior. It allows a larger range of step-sizes to be used, and creates its own oscillations. What is going on?

## [Why Momentum Really Works](#)



Step-size  $\alpha = 0.0030$

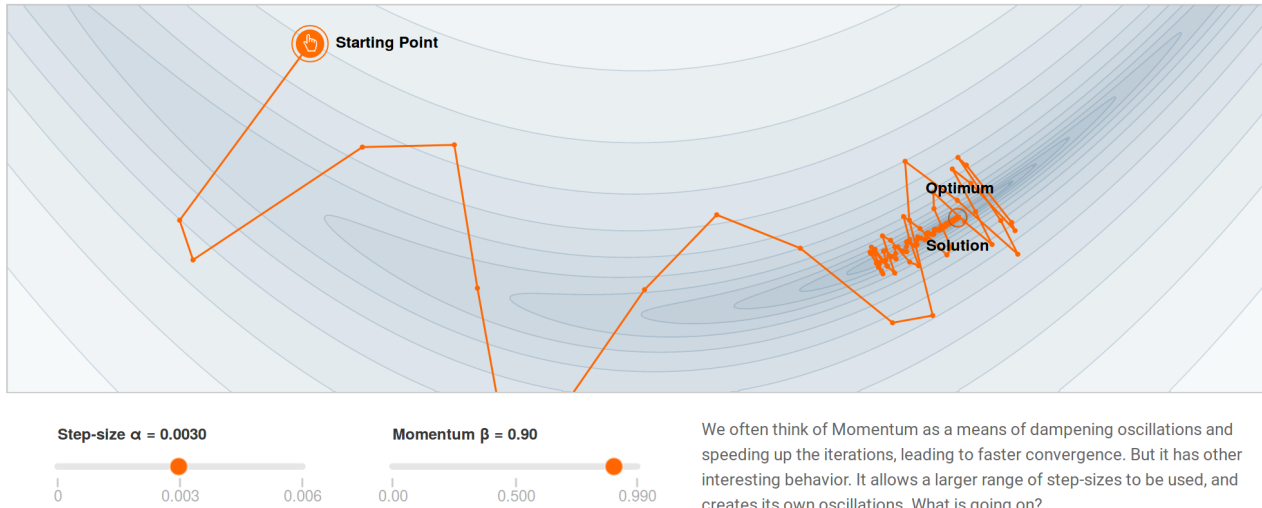


Momentum  $\beta = 0.80$



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



## Why Momentum Really Works

# Alternative optimizers

- SGD (with Nesterov momentum)
  - Simple to implement
  - Very sensitive to initial value of  $\eta$
  - Need learning rate scheduling

# Alternative optimizers

- SGD (with Nesterov momentum)
  - Simple to implement
  - Very sensitive to initial value of  $\eta$
  - Need learning rate scheduling
- Adam: adaptive learning rate scale for each param
  - Global  $\eta$  set to  $3e-4$  often works well enough
  - Good default choice of optimizer (often)

# Alternative optimizers

- SGD (with Nesterov momentum)
  - Simple to implement
  - Very sensitive to initial value of  $\eta$
  - Need learning rate scheduling
- Adam: adaptive learning rate scale for each param
  - Global  $\eta$  set to  $3e-4$  often works well enough
  - Good default choice of optimizer (often)
- But well-tuned SGD with LR scheduling can generalize better than Adam (with naive l2 reg)...

# Alternative optimizers

- SGD (with Nesterov momentum)
  - Simple to implement
  - Very sensitive to initial value of  $\eta$
  - Need learning rate scheduling
- Adam: adaptive learning rate scale for each param
  - Global  $\eta$  set to  $3e-4$  often works well enough
  - Good default choice of optimizer (often)
- But well-tuned SGD with LR scheduling can generalize better than Adam (with naive l2 reg)...
- Promising stochastic second order methods: [K-FAC](#) and [Shampoo](#) can be used to accelerate training of very large models.

# The Karpathy Constant for Adam



**Andrej Karpathy** 

@karpathy

Following



3e-4 is the best learning rate for Adam, hands down.

4:01 AM - 24 Nov 2016

101 Retweets 408 Likes



23



101

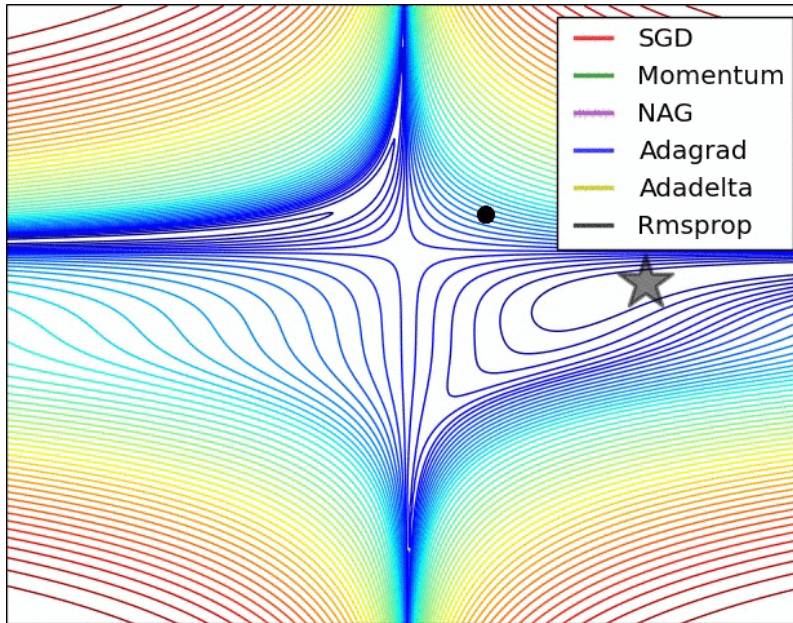


408





# Optimizers around a saddle point



Credits: Alec Radford