# Intro to ML

December 1st, 2021

# CHAPTER 11:
# Multilayer Perceptrons

# Learning Boolean AND

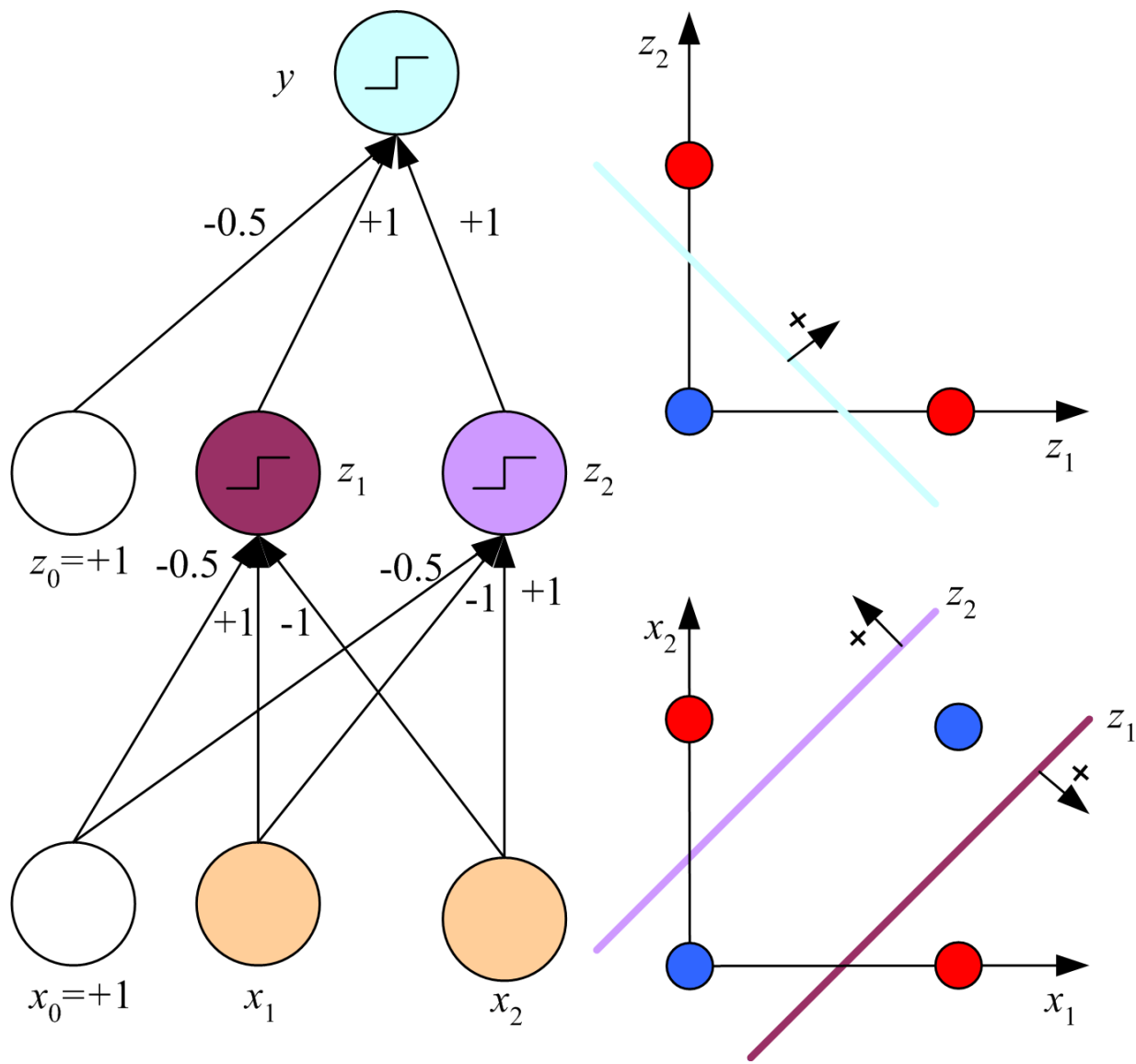| $x_1$ | $x_2$ | $r$ |
|-------|-------|-----|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

$y$

$-1.5$

$+1$    $+1$

$x_0 = +1$    $x_1$    $x_2$

$x_2$

$1.5$

$(0,1)$    $(1,1)$

$+$

$(0,0)$    $(1,0)$    $1.5$    $x_1$

$y$

$-0.5$  $+1$  $+1$

$z_0=+1$  $-0.5$  $z_1$  $-0.5$  $z_2$

$+1$  $-1$  $-1$  $+1$

$x_0=+1$  $x_1$  $x_2$

$z_2$  $z_1$

$x_2$  $z_2$  $z_1$  $x_1$

| $x_1$ | $x_2$ | $z_1$ | $z_2$ | $y$ |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 | 1 |
| 1 | 0 | 1 | 0 | 1 |
| 1 | 1 | 0 | 0 | 0 |

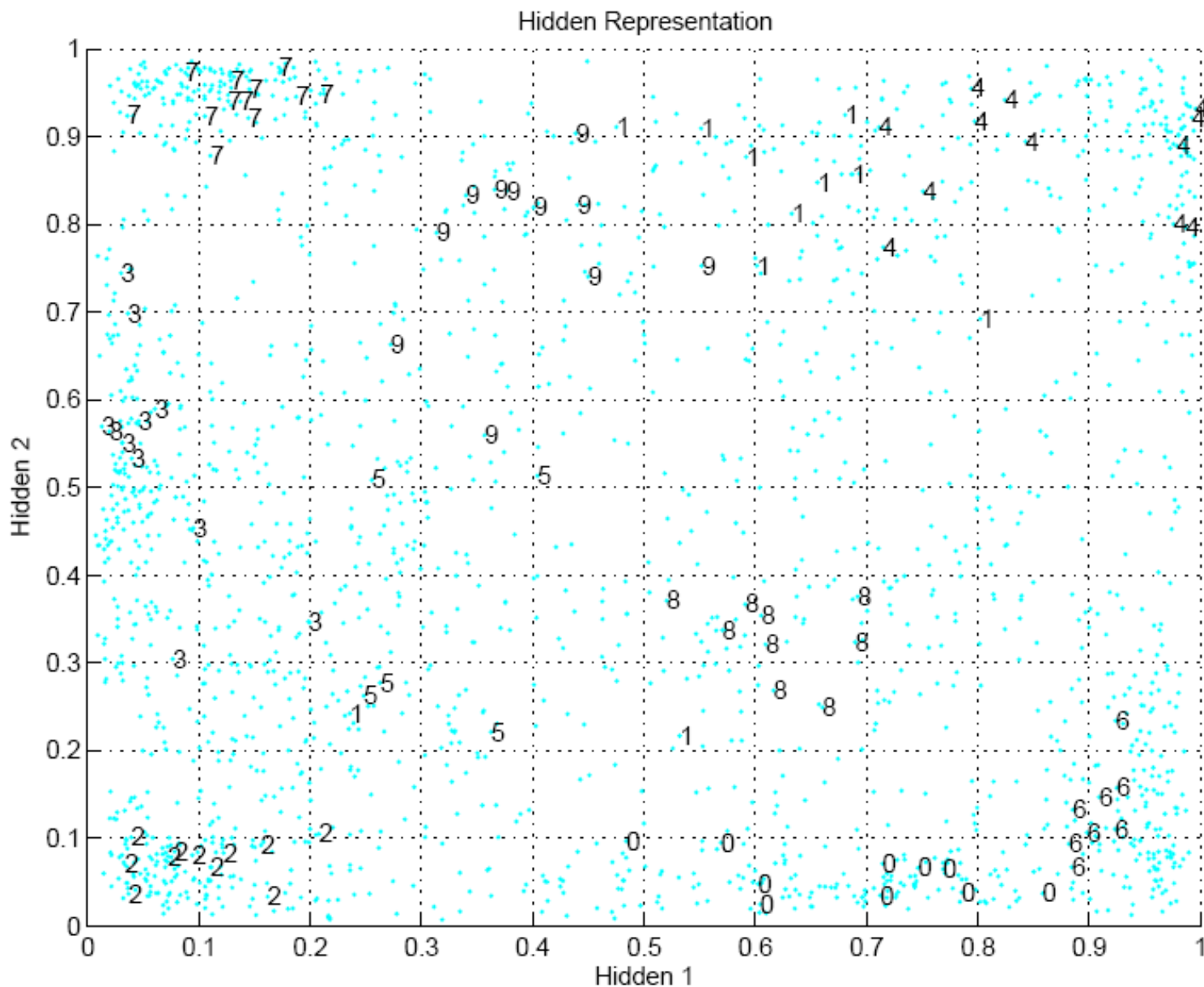$x_1$ XOR $x_2$ = ($x_1$ AND ~$x_2$) OR (~$x_1$ AND $x_2$)

4

# Learning Hidden Representations

- MLP is a generalized linear model where hidden units are the nonlinear basis functions:

$$y = \sum_{h=1}^{H} v_h \phi(\boldsymbol{x}|\boldsymbol{w}_h)$$

where

$$\phi(\boldsymbol{x}|\boldsymbol{w}_h) \equiv \text{sigmoid}(\boldsymbol{w}_h^T \boldsymbol{x})$$

- The advantage is that the basis function parameters can also be learned from data.

- The hidden units, $z_h$, learn a *code/embedding*, a representation in the hidden space
  - If H space is < original d dimension: dimension reduction

- Transfer learning: Use code in another task

- Semisupervised learning: Transfer from an unsupervised to supervised problem
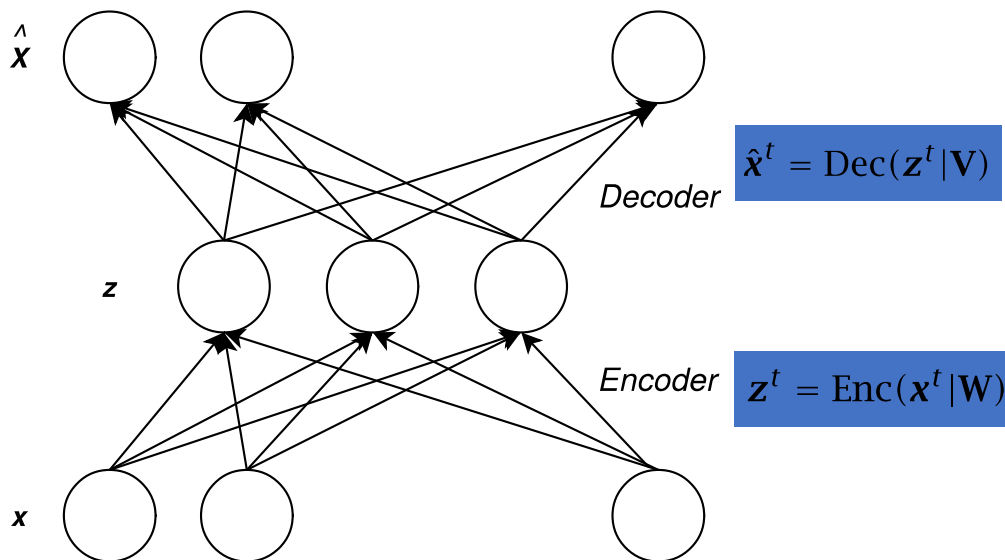
100 data points

64 inputs
2 hidden units

10 output

# Autoencoders

An MLP structure

Equal number of input and ouput
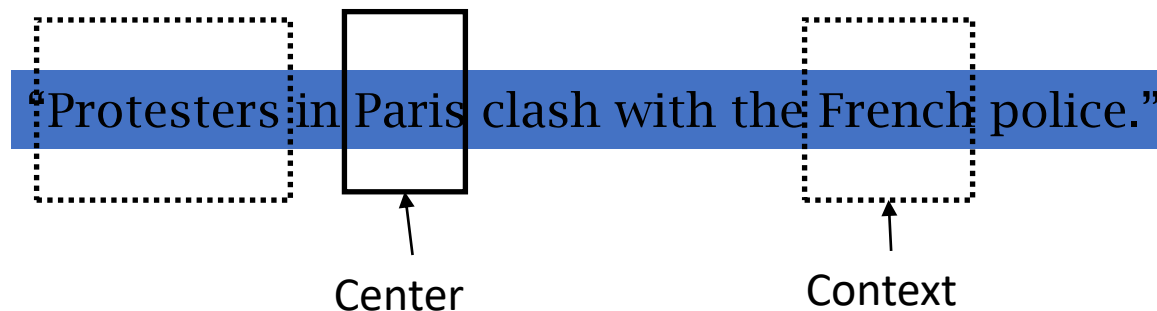
Dimension reduction



$\hat{\boldsymbol{x}}^t = \mathrm{Dec}(\boldsymbol{z}^t | \mathbf{V})$

*Decoder*

$\boldsymbol{z}^t = \mathrm{Enc}(\boldsymbol{x}^t | \mathbf{W})$

*Encoder*

$$E(\mathbf{W}, \mathbf{V} | \mathcal{X}) = \sum_t \| \boldsymbol{x}^t - \hat{\boldsymbol{x}}^t \|^2 = \sum_t \| \boldsymbol{x}^t - \mathrm{Dec}(\mathrm{Enc}(\boldsymbol{x}^t | \mathbf{W}) | \mathbf{V}) \|^2$$

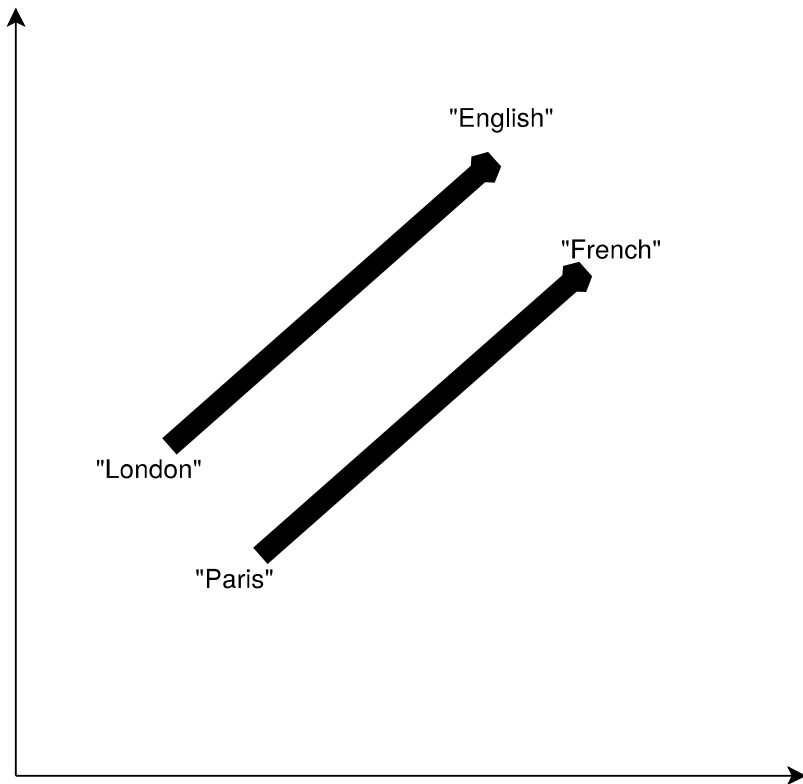- Variants: Denoising, sparse autoencoders

# Example: Word2vec

- Learn an embedding for words for NLP
- Skipgram: An autencoder with linear encoder and decoder where input is the center word and output is a context word

"Protesters in Paris clash with the French police."

Center

Context

- Similar words appear in similar contexts, so similar codes will be learned for them

# Vector Algebra



Because they will always appear in similar contexts in pairs in a large corpus, we expect

vec("English")-vec("London")

to be similar to

vec("French")-vec("Paris")

so,

vec("Paris") + vec("English")-vec("London")

will be similar to

vec("French")

CHAPTER 12:

# Deep Learning

# Deep Neural Networks
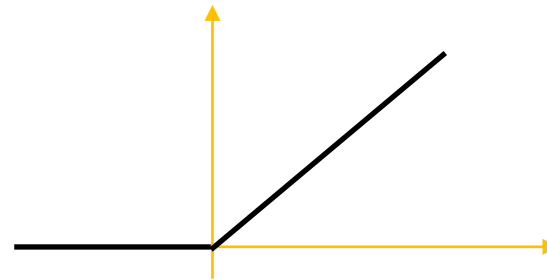
output

$\uparrow$

$\vdots$

input

- Many hidden layers
  - Problematic in training: chain rule multiplication of derivatives (vanish of gradients or explosion)
- End-to-end training
- Learn increasingly abstract representations with minimal human contribution
  - Layers of abstraction in intuitive
  - Vision, speech, language, and so on

Representation learning is really the KEY behind Deep learning

# Rectified Linear Unit (ReLU)

Left deriviation:
Relu'(a) = 1 if a>0, 0 otherwise

$$\text{ReLU}(a) = \begin{cases} a & \text{if } a > 0 \\ 0 & \text{otherwise} \end{cases}$$
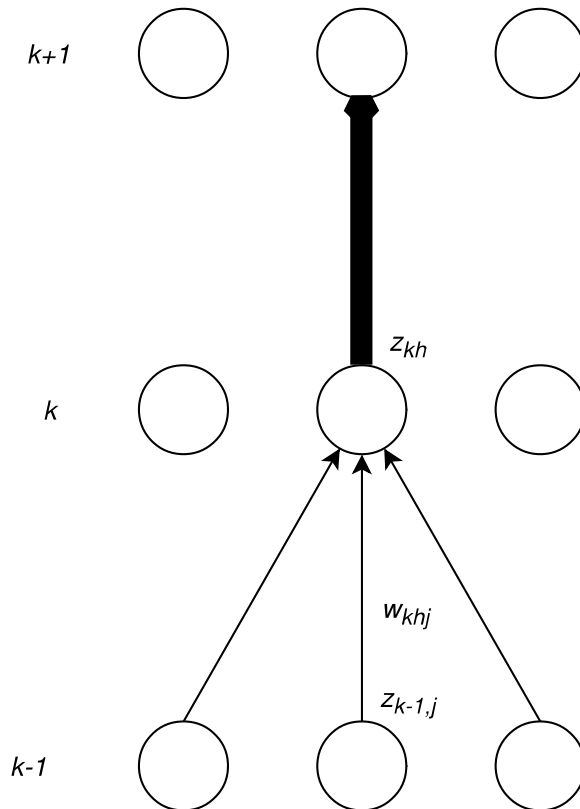
- Does not saturate for large *a*

- Leads to a sparse representation
  - Some of the nodes will results in 0

- No learning for *a*<0, be careful with initialization
  - Initialization should make sure all weights are positive

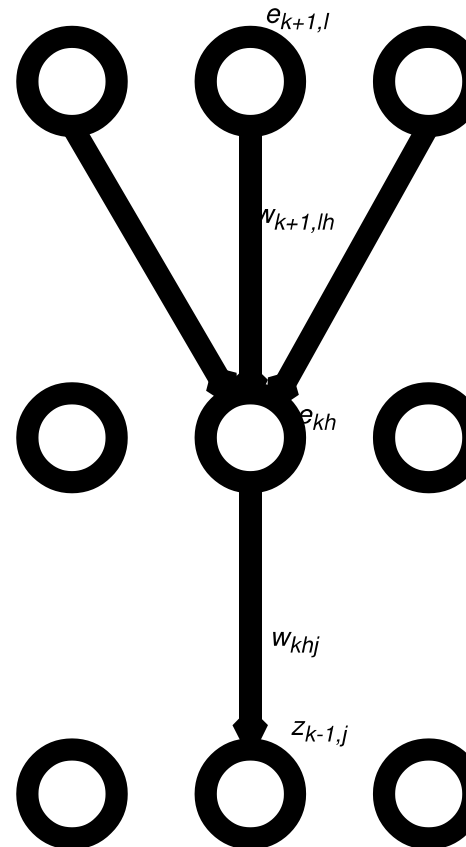- Leaky ReLU: {a if a>0, $\alpha$a otherwise, $\alpha$ is set 0.01)

Some derivative left

# Generalizing Backpropagation

$$z_{kh}^t = f\left(\sum_j w_{khj} z_{k-1,j}^t\right)$$

$$e_{kh}^t = \sum_l e_{k+1,l}^t w_{k+1,lh}$$



$k+1$

$e_{k+1,l}$

$w_{k+1,lh}$

$z_{kh}$

$k$

$e_{kh}$

$w_{khj}$

$w_{khj}$

$z_{k-1,j}$

$z_{k-1,j}$

$k-1$

Forward

Backward

13

$$\Delta w_{khj} = \eta \sum_t e^t_{kh} f'(z^t_{kh}) z^t_{k-1,j}$$

- If k is the output layer, there is no layer k+1
  - $e^t_{kh}$ is the derivative of the error function with respect to the output $z_{kh}$.

- If not, follow the previous backward rules
- f' is the derivative of the activation function

# Improving Training Convergence

- Gradient decent. Simply, local, changes uses only the values of the presynaptic and postsynaptic and the error (suitably backpropagated) – can converge slowly

- Momentum. At each update, also add a fraction of the average of past gradients (t denotes iteration):

$$s_i^t = \alpha s_i^{t-1} + (1 - \alpha)\frac{\partial E^t}{\partial w_i}$$

Running average of past gradient

$$\Delta w_i^t = -\eta s_i^t$$

Update using such an average

# Adaptive Learning Factor

- RMSprop. Make update inversely proportional to <u>the sum of past gradients</u>, so, update more when gradient is small and less where it is large.

$$\Delta w_i^t = -\frac{\eta}{\sqrt{r_i^t}} \frac{\partial E^t}{\partial w_i}$$

where $r_i$ is the accumulated past gradient,

$$r_i^t = \rho r_i^{t-1} + (1 - \rho) \left| \frac{\partial E^t}{\partial w_i} \right|^2$$

p is generally set at 0.999

# ADAM: Adaptive Learning Factor w/ Momentum

First moment – just the gradient

$$s_i^t \quad = \quad \alpha s_i^{t-1} + (1-\alpha)\frac{\partial E^t}{\partial w_i}$$

$$\Delta w_i^t \quad = \quad -\eta \frac{\tilde{s}_i^t}{\sqrt{\tilde{r}_i^t}}$$

$$r_i^t \quad = \quad \rho r_i^{t-1} + (1-\rho)\left|\frac{\partial E^t}{\partial w_i}\right|^2$$

Second moment – square of gradient

Initially both $s$ and $r$ terms are 0, so we bias-correct them (both $\alpha$ and $\rho$ are <1, so they get smaller as $t$ gets large):

$$\tilde{s}_i^t \quad = \quad \frac{s_i^t}{1-\alpha^t} \text{ and } \tilde{r}_i^t = \frac{r_i^t}{1-\rho^t}$$