# Lecture 22: Neural Networks, Back-propagation, etc

Instructor: Prof. Ganesh Ramakrishnan

inputs

$1$

$x_1$

$x_2$

$x_n$

$w_{01}$

$w_{11}$

$w_{21}$

$w_{n1}$

$w_{02}$

$w_{12}$

$w_{22}$

$w_{n2}$

$z_1 = g\left(\sum\right)$

$z_2 = g\left(\sum\right)$

$f_1 = g(.)$

$f_2 = g(.)$

One hot encoding $\Rightarrow$ Exactly one of the K nodes in $(L^{th})$ o/p layer expected to be a 1 & rest 0

STEP 0: Pick a network architecture

- Number of input units: Dimension of features $\phi\left(\mathbf{x}^{(i)}\right)$.
- Number of output units: Number of classes.
- Reasonable default: 1 hidden layer, or if $> 1$ hidden layer, have same number of hidden units in every layer.
- Number of hidden units in each layer a constant factor (3 or 4) of dimension of $x$.
- We will use
  - the smooth sigmoidal function $g(s) = \frac{1}{1+e^{-s}}$: **We have now learnt how to train a single node sigmoidal (LR) neural network**
  - instead of the non-smooth step function $g(s) = 1$ if $s \in [\theta, \infty)$ and $g(s) = 0$ otherwise.
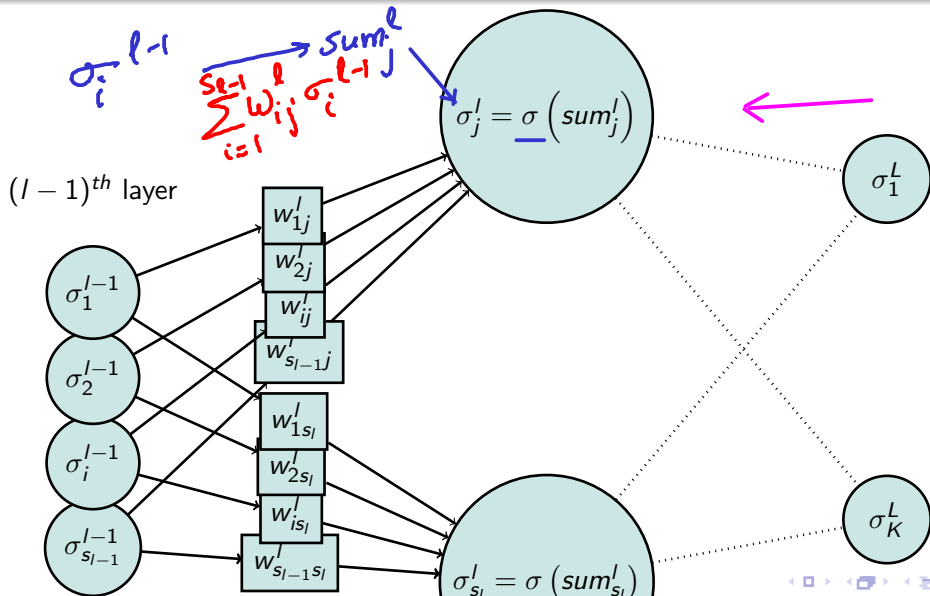
1. Randomly initialize weights $w_{ij}^l$ for $l = 1, \ldots, L$, $i = 1, \ldots, s_l$, $j = 1, \ldots, s_{l+1}$.
2. Implement **forward propagation** to get $f_w(\mathbf{x})$ for any $x \in \mathcal{D}$.
3. Execute **backpropagation**
   1. by computing partial derivatives $\frac{\partial}{\partial w_{ij}^{(l)}} E(w)$ for $l = 1, \ldots, L$, $i = 1, \ldots, s_l$, $j = 1, \ldots, s_{l+1}$.
   2. and using gradient descent to try to minimize (non-convex) $E(w)$ as a function of parameters $\mathbf{w}$.

$$w_{ij}^l = w_{ij}^l - \eta \frac{\partial}{\partial w_{ij}^{(l)}} E(w)$$

4. Verify that the cost function $E(w)$ has indeed reduced, else resort to some random perturbation of weights $\mathbf{w}$.

- The Neural Network objective to be minimized:

$$E(\mathbf{w}) = -\frac{1}{m} \left[ \sum_{i=1}^{m} \sum_{k=1}^{K} y_k^{(i)} \log\left(\sigma_k^L\left(\mathbf{x}^{(i)}\right)\right) + \left(1 - y_k^{(i)}\right) \log\left(1 - \sigma_k^L\left(\mathbf{x}^{(i)}\right)\right) \right] \tag{1}$$

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L} \sum_{i=1}^{s_{l-1}} \sum_{j=1}^{s_l} \left(w_{ij}^l\right)^2$$

Average cross entropy at $L^{th}$ layer

Regularization on wts across all layers

# Gradient Computation

- The Neural Network objective to be minimized:

$$E(\mathbf{w}) = -\frac{1}{m}\left[\sum_{i=1}^{m}\sum_{k=1}^{K} y_k^{(i)} \log\left(\sigma_k^L\left(\mathbf{x}^{(i)}\right)\right) + \left(1 - y_k^{(i)}\right)\log\left(1 - \sigma_k^L\left(\mathbf{x}^{(i)}\right)\right)\right]$$

$$+\frac{\lambda}{2m}\sum_{l=1}^{L}\sum_{i=1}^{s_{l-1}}\sum_{j=1}^{s_l}\left(w_{ij}^l\right)^2 \tag{1}$$

- $sum_j^l = \sum_{k=1}^{s_{l-1}} w_{kj}^l \sigma_k^{l-1}$ and $\sigma_i^l = \frac{1}{1 + e^{-sum_i^l}}$

- $\frac{\partial \mathbf{E}}{\partial \mathbf{w_{ij}^l}} = \frac{\partial \mathbf{E}}{\partial \sigma_j^l} \frac{\partial \sigma_j^l}{\partial \mathbf{sum_j^l}} \frac{\partial \mathbf{sum_j^l}}{\partial \mathbf{w_{ij}^l}} + \frac{\lambda}{2m} w_{ij}^l$

- $\frac{\partial \sigma_j^l}{\partial \mathbf{sum_j^l}} = \left(\frac{1}{1 + e^{-sum_i^l}}\right)\left(1 - \frac{1}{1 + e^{-sum_i^l}}\right) = \sigma_j^l (1 - \sigma_j^l)$

- $\frac{\partial \mathbf{sum_j^l}}{\partial \mathbf{w_{ij}^l}} = \frac{\partial}{\partial w_{ij}^l}\left(\sum_{k=1}^{s_{l-1}} w_{kj}^l \sigma_k^{l-1}\right) = \sigma_i^{l-1}$

- For a single example $(\mathbf{x}, y)$:

$$-\left[\sum_{k=1}^{K} y_k \log\left(\sigma_k^L(\mathbf{x})\right) + (1-y_k)\log\left(1-\sigma_k^L(\mathbf{x})\right)\right]$$

$$+\frac{\lambda}{2m}\sum_{l=1}^{L}\sum_{i=1}^{s_{l-1}}\sum_{j=1}^{s_l}\left(w_{ij}^l\right)^2 \qquad (2)$$

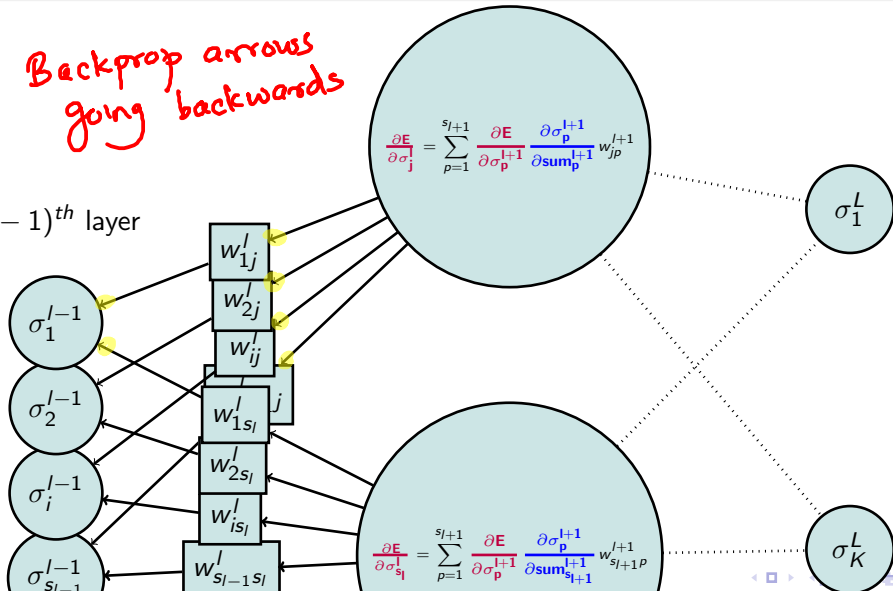*(handwritten: σ of current layer)*

*(handwritten: sum of current layers)*

- $\frac{\partial \mathbf{E}}{\partial \sigma_j^l} = \sum_{p=1}^{s_{l+1}} \frac{\partial E}{\partial sum_p^{l+1}} \frac{\partial sum_p^{l+1}}{\partial \sigma_j^l} = \sum_{p=1}^{s_{l+1}} \frac{\partial \mathbf{E}}{\partial \sigma_p^{l+1}} \frac{\partial \sigma_p^{l+1}}{\partial \mathbf{sum_p^{l+1}}} w_{jp}^{l+1}$ since $\frac{\partial sum_p^{l+1}}{\partial \sigma_j^l} = w_{jp}^{l+1}$

- $\frac{\partial \mathbf{E}}{\partial \sigma_j^L} = -\frac{y_j}{\sigma_j^L} - \frac{1-y_j}{1-\sigma_j^L}$

*(handwritten: σ of prev layer)*

Backprop arrows going backwards

$(l-1)^{th}$ layer

$$\frac{\partial \mathbf{E}}{\partial \sigma_{\mathbf{j}}^{\mathbf{l}}} = \sum_{p=1}^{s_{l+1}} \frac{\partial \mathbf{E}}{\partial \sigma_{\mathbf{p}}^{\mathbf{l+1}}} \frac{\partial \sigma_{\mathbf{p}}^{\mathbf{l+1}}}{\partial \mathbf{sum}_{\mathbf{p}}^{\mathbf{l+1}}} w_{jp}^{l+1}$$

$$\frac{\partial \mathbf{E}}{\partial \sigma_{\mathbf{s_l}}^{\mathbf{l}}} = \sum_{p=1}^{s_{l+1}} \frac{\partial \mathbf{E}}{\partial \sigma_{\mathbf{p}}^{\mathbf{l+1}}} \frac{\partial \sigma_{\mathbf{p}}^{\mathbf{l+1}}}{\partial \mathbf{sum}_{\mathbf{s_{l+1}}}^{\mathbf{l+1}}} w_{s_{l+1}p}^{l+1}$$

$\sigma_1^{l-1}$

$\sigma_2^{l-1}$

$\sigma_i^{l-1}$

$\sigma_{s_{l-1}}^{l-1}$

$w_{1j}^l$

$w_{2j}^l$

$w_{ij}^l$

$w_{1s_l}^l$

$w_{2s_l}^l$

$w_{is_l}^l$

$w_{s_{l-1}s_l}^l$

$j$

$\sigma_1^L$

$\sigma_K^L$

$(l-1)^{th}$ layer

$$\frac{\partial \mathbf{E}}{\partial \mathbf{w_{ij}^l}} = \frac{\partial \mathbf{E}}{\partial \sigma_j^l} \frac{\partial \sigma_j^l}{\partial \mathbf{sum_j^l}} \frac{\partial \mathbf{sum_j^l}}{\partial \mathbf{w_{ij}^l}} + \frac{\lambda}{2m} w_{ij}^l$$

$\frac{\partial \mathbf{E}}{\partial \sigma_j^l}$

$\frac{\partial \mathbf{E}}{\partial \sigma_{s_l}^l}$
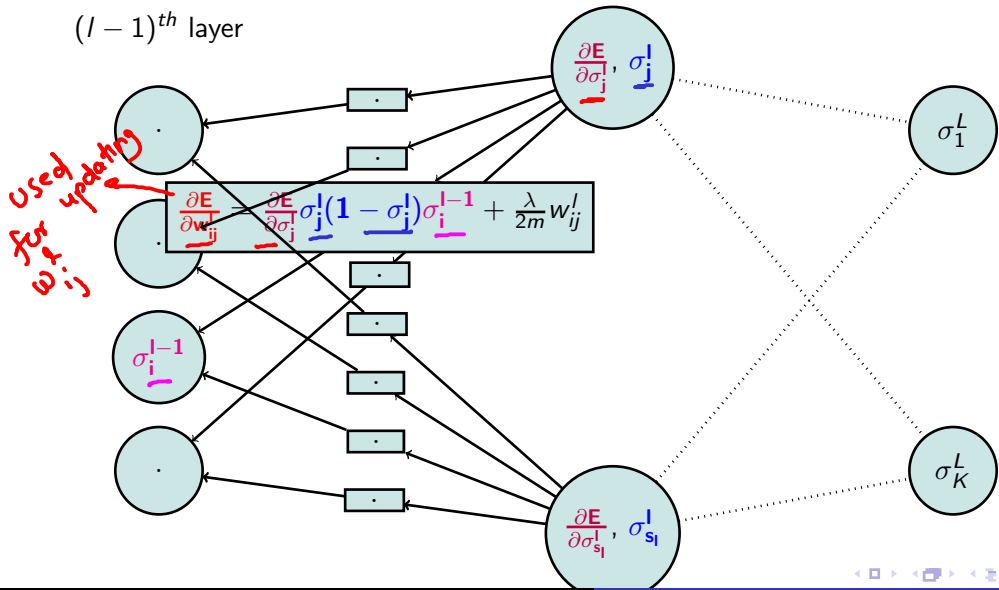
$\sigma_1^L$

$\sigma_K^L$

$\sigma_i^{l-1}$

# Recall and Substitute

- $sum_j^l = \sum_{k=1}^{s_{l-1}} w_{kj}^l \sigma_k^{l-1}$ and $\sigma_i^l = \frac{1}{1+e^{-sum_i^l}}$

- $\frac{\partial \mathbf{E}}{\partial \mathbf{w_{ij}^l}} = \frac{\partial \mathbf{E}}{\partial \sigma_j^l} \frac{\partial \sigma_j^l}{\partial \mathbf{sum_j^l}} \frac{\partial \mathbf{sum_j^l}}{\partial \mathbf{w_{ij}^l}} + \frac{\lambda}{2m} w_{ij}^l$

- $\frac{\partial \sigma_j^l}{\partial \mathbf{sum_j^l}} = \sigma_j^l (\mathbf{1} - \sigma_j^l)$

- $\frac{\partial \mathbf{sum_j^l}}{\partial \mathbf{w_{ij}^l}} = \sigma_i^{l-1}$

- $\frac{\partial \mathbf{E}}{\partial \sigma_j^l} = \sum_{p=1}^{s_{l+1}} \frac{\partial \mathbf{E}}{\partial \sigma_j^{l+1}} \frac{\partial \sigma_j^{l+1}}{\partial \mathbf{sum_j^{l+1}}} w_{jp}^{l+1}$

- $\frac{\partial \mathbf{E}}{\partial \sigma_j^{\mathbf{L}}} = -\frac{\mathbf{y_j}}{\sigma_j^{\mathbf{L}}} - \frac{\mathbf{1-y_j}}{\mathbf{1}-\sigma_j^{\mathbf{L}}}$

$(l-1)^{th}$ layer

Used for updating $w_{ij}$

$$\frac{\partial \mathbf{E}}{\partial \mathbf{w_{ij}^l}} = \frac{\partial \mathbf{E}}{\partial \sigma_j^l} \sigma_j^l (1 - \sigma_j^l) \sigma_i^{l-1} + \frac{\lambda}{2m} w_{ij}^l$$

$\frac{\partial \mathbf{E}}{\partial \sigma_j^l}, \sigma_j^l$

$\frac{\partial \mathbf{E}}{\partial \sigma_{s_l}^l}, \sigma_{s_l}^l$

$\sigma_i^{l-1}$

$\sigma_1^L$

$\sigma_K^L$

$(l-1)^{th}$ layer

$\dfrac{\partial \mathbf{E}}{\partial \sigma_j^l}$

$w_{ij}^l = w_{ij}^l - \eta \dfrac{\partial \mathbf{E}}{\partial \mathbf{w}_{ij}^l}$

$\sigma_i^{l-1}$

$\dfrac{\partial \mathbf{E}}{\partial \sigma_{s_l}^l}$

$\sigma_1^L$

$\sigma_K^L$

Forward inference (before beginning backprop)

① The values $w_{ij}^{l-1} .. w_{ij}^l$ are (implicitly) used from previous layers & previous iterations

Back prop ② The partial derivatives $\dfrac{\partial \sigma_j^l}{\partial sum_j^l}$, $\dfrac{\partial E}{\partial \sigma_j^l}$ are used from current iterations but for $l' > l$ (following layers)

All $w_{ij}^l$ for fixed $l$ updated simultaneously

$W_{ij}^l$'s being updated simultaneously etc and other addition/multiplication operations are converted into efficient matrix/tensor operations in NN packages such as tensorflow, Theano, Torch etc.

# The Backpropagation Algorithm for Training NN

1. Randomly initialize weights $w_{ij}^l$ for $l = 1, \ldots, L$, $i = 1, \ldots, s_l$, $j = 1, \ldots, s_{l+1}$. *→ Typically 0's, unless RELU is used*

2. Implement **forward propagation** to get $f_{\mathbf{w}}(\mathbf{x})$ for every $\mathbf{x} \in \mathcal{D}$. *(using $w_{ij}^l$ from previous iterations)*

3. Execute **backpropagation** on any misclassified $\mathbf{x} \in \mathcal{D}$ by performing gradient descent to minimize (non-convex) $E(\mathbf{w})$ as a function of parameters $\mathbf{w}$.

4. $\frac{\partial \mathbf{E}}{\partial \sigma_j^L} = -\frac{\mathbf{y_j}}{\sigma_j^L} - \frac{1-\mathbf{y_j}}{1-\sigma_j^L}$ for $j = 1$ to $s_L$. *(on final/output layer)*

5. For $l = L - 1$ down to 2:

   *Illustrated on prev slides*

   1. $\frac{\partial \mathbf{E}}{\partial \sigma_j^l} = \sum_{p=1}^{s_{l+1}} \frac{\partial \mathbf{E}}{\partial \sigma_j^{l+1}} \sigma_j^{l+1}(\mathbf{1} - \sigma_j^{l+1}) w_{jp}^{l+1}$

   2. $\frac{\partial \mathbf{E}}{\partial \mathbf{w_{ij}^l}} = \frac{\partial \mathbf{E}}{\partial \sigma_j^l} \sigma_j^l(\mathbf{1} - \sigma_j^l)\sigma_i^{l-1} + \frac{\lambda}{2m} w_{ij}^l$

   3. $w_{ij}^l = w_{ij}^l - \eta \frac{\partial \mathbf{E}}{\partial \mathbf{w_{ij}^l}}$

   *Backpropagate on previous layers successively*

6. Keep picking misclassified examples until the cost function $E(\mathbf{w})$ shows significant reduction; else resort to some random perturbation of weights $\mathbf{w}$ and restart a couple of times.
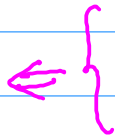
Note: If $\sigma_j^\ell(sum_j^\ell) = \dfrac{1}{1 + e^{-sum_j^\ell}}$ is

replaced by RELU (Reetified linear unit)

$$relu_j^\ell(sum_j^\ell) = max(0, sum_j^\ell)$$

Then $\dfrac{\partial \sigma_j^\ell}{\partial sum_j^\ell}$ becomes $\dfrac{\partial relu_j^\ell}{\partial sum_j^\ell} = 1$ (if $sum_j^\ell > 0$)

Therefore initializing $w_{ij}^\ell = 0$ can become meaningless. $\Longleftarrow \Big\{$

subgradient

$\& \dfrac{\partial relu_j^\ell}{\partial sum_j^\ell} = 0$ (if $sum_j^\ell \leq 0$)

SGD on (Strongly) Convex problems $\Rightarrow$ Pending error is $(O(\frac{1}{k}))$ $O(\frac{1}{\sqrt{k}})$ $\longrightarrow$ $k =$ #iterations

- **Average (sigmoidal) gradient on a minibatch of $m_b$ examples:**

$$\frac{\partial \mathbf{E}}{\partial \sigma_j^L} = \frac{1}{m_b} \sum_{i=1}^{m_b} -\frac{y_j^{(i)}}{\sigma_j^L(x^{(i)})} - \frac{1 - y_j^{(i)}}{1 - \sigma_j^L(x^{(i)})} \text{ for } j = 1 \text{ to } s_L.$$

only indicative
Does not say much
about our nonconvex
avg cross entropy
fn for NN's
with hidden layers.

$\nabla E(i \in \{m\})$   $\nabla E(i \in \{m_b\}) \cdots \nabla E(i)$

True or batch   Minibatch Stochastic   Stochastic

Necessary
condition for optimality : $\nabla E(i \in \{m\}) = 0$

(Does not imply) $\Rightarrow$   $\nabla E(i \in \{m_b\}) = 0$ OR $\nabla E(i) = 0$

**More probably 0 than for stochastic**

$$\omega^{(k)} = \omega^{(k-1)} + \eta \nabla E(m_b)$$

If $m_b = m$, then as $\omega^{(k-1)} \to \omega^*$

$$\nabla_m E(\omega^{(k-1)}) \to 0$$

For $m_b < m$, can

$\eta_k$ compensate for

a misbehaving

$\nabla E_{m_b}$

① $\eta_k$ instead of $\eta$

② $\eta_k \to 0$ as $k\uparrow$

$$\Rightarrow \omega^{(k-1)} \to \omega^{(k)} \to \omega^* \text{ as } k\uparrow$$

If $m_b < m$, then even as $\omega^{(k-1)} \to \omega^*$

No guarantee that $\nabla_{m_b} E(\omega^{(k-1)}) \to 0$

$\Rightarrow$ Even if $\omega^{(k-1)} \to \omega^{(*)}$ yet $\omega^{(k)}$ is updated

to drift away from $\omega^*$

# Efficient alternatives to Stochastic Gradient Descent (SGD)

SGD on (Strongly) Convex problems $\Rightarrow$ Pending error is $(O(\frac{1}{k}))$ $O(\frac{1}{\sqrt{k}})$

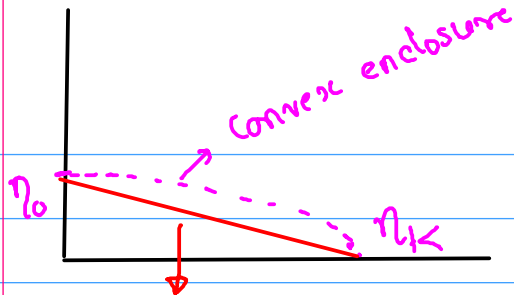- **Average (sigmoidal) gradient on a minibatch of $m_b$ examples:**

$$\frac{\partial \mathbf{E}}{\partial \sigma_j^L} = \frac{1}{m_b} \sum_{i=1}^{m_b} -\frac{\mathbf{y}_j^{(i)}}{\sigma_j^L(\mathbf{x}^{(i)})} - \frac{1 - \mathbf{y}_j^{(i)}}{1 - \sigma_j^L(\mathbf{x}^{(i)})} \text{ for } j = 1 \text{ to } s_L.$$

- **Adaptive Learning Rate:** Expect at optimality $\nabla \mathbf{E} = 0$. But not if gradient is approximated on a sample $m_b$! A sufficient condition for convergence of SGD is:
  - $\eta_k$ (learning rate) vary across iterations $w_{ij}^l = w_{ij}^l - \eta_k \frac{\partial \mathbf{E}}{\partial \mathbf{w}_{ij}^l}$
  - $\sum_{k=1}^{\infty} \eta_k = \infty$ and $\sum_{k=1}^{\infty} \eta_k^2 \leq \gamma < \infty$ (eg: $\eta_k = \frac{1}{k}$, $\gamma = \frac{\pi^2}{6}$). Commonly $\eta_k = (1 - \beta_k)\eta_0 + \beta_k \eta_K$, with $\beta_k = \frac{k}{K}$ that is convex combination of some max ($\eta_0$) and min ($\eta_K$) values.

- **Adagrad:** Individually adapts learning rates of all model parameters by scaling them inversely proportional to square root of sum of all their historical squared values $\eta_k = 1/\sqrt{\|\nabla E\|^2}$ over prev iterations averaged

- **RMSProp:** Modifies AdaGrad for non-convex setting by changing the gradient accumulation into an exponentially weighted moving average

Convex enclosure

$\eta_0$

$\eta_k$

Line segment for $\eta_k$

$$\eta_k = \beta \eta_0 + (1-\beta) \eta_K$$

$\dfrac{k}{K}$

$0.5$

Assume $K$ is large enough & stopping criterion is that error on validation has stopped decreasing any further

*Adagrad: Momentum through learning rate*

- **<u>Momentum[1] Accelerated SGD</u>:** $w_{ij}^l = w_{ij}^l + v_{ij}^l$ where $v_{ij}^l = \alpha v_{ij}^l - \eta_k \frac{\partial \mathbf{E}}{\partial w_{ij}^l}$

  *Directly capture momentum for $w_{ij}$*

- **Different choice of Activation Function $\sigma$ such as ReLU, tanh, etc**

- **<u>Adam</u>:** Best seen as a variant on the combination of RMSProp and Momentum

The following part remains the same with sigmoidal functions (for $l = L - 1$ down to 2):

① $\frac{\partial \mathbf{E}}{\partial \sigma_j^l} = \sum_{p=1}^{s_{l+1}} \frac{\partial \mathbf{E}}{\partial \sigma_j^{l+1}} \sigma_j^{l+1}(1 - \sigma_j^{l+1}) w_{jp}^{l+1}$

② $\frac{\partial \mathbf{E}}{\partial w_{ij}^l} = \frac{\partial \mathbf{E}}{\partial \sigma_j^l} \sigma_j^l (1 - \sigma_j^l) \sigma_i^{l-1} + \frac{\lambda}{2m} w_{ij}^l$

*Associated with $\eta_k$.*   *associated with $w_{ij}$*

---

[1]In practice, $\alpha \in \{0.5, 0.9, 0.99\}$