NEURONS

## GD / SGD For Linear Neuron (Linear activation)

| GD | SGD |
|---|---|
| $(\boldsymbol{X}, \boldsymbol{d})$ | $(\boldsymbol{x}_p, d_p)$ |
| $J = \dfrac{1}{2}\sum_{p=1}^{P}(d_p - y_p)^2$ | $J = \dfrac{1}{2}(d_p - y_p)^2$ |
| $\boldsymbol{y} = \boldsymbol{u} = \boldsymbol{X}\boldsymbol{w} + b\boldsymbol{1}_P$ | $y_p = u_p = \boldsymbol{x}_p^{T}\boldsymbol{w} + b$ |
| $\boldsymbol{w} \leftarrow \boldsymbol{w} + \alpha\boldsymbol{X}^{T}(\boldsymbol{d} - \boldsymbol{y})$ | $\boldsymbol{w} \leftarrow \boldsymbol{w} + \alpha(d_p - y_p)\boldsymbol{x}_p$ |
| $b \leftarrow b + \alpha\boldsymbol{1}_P^{T}(\boldsymbol{d} - \boldsymbol{y})$ | $b \leftarrow b + \alpha(d_p - y_p)$ |

## GD / SGD For Perceptron (Sigmoid activation)

| GD | SGD |
|---|---|
| $(\boldsymbol{X}, \boldsymbol{d})$ | $(\boldsymbol{x}_p, d_p)$ |
| $J = \dfrac{1}{2}\sum_{p=1}^{P}(d_p - y_p)^2$ | $J = \dfrac{1}{2}(d_p - y_p)^2$ |
| $\boldsymbol{u} = \boldsymbol{X}\boldsymbol{w} + b\boldsymbol{1}_P$ | $u_p = \boldsymbol{x}_p^{T}\boldsymbol{w} + b$ |
| $\boldsymbol{y} = f(\boldsymbol{u})$ | $y_p = f(u_p)$ |
| $\boldsymbol{w} = \boldsymbol{w} + \alpha\boldsymbol{X}^{T}(\boldsymbol{d} - \boldsymbol{y})\cdot f'(\boldsymbol{u})$ | $\boldsymbol{w} = \boldsymbol{w} + \alpha(d_p - y_p)f'(u_p)\boldsymbol{x}_p$ |
| $b = b + \alpha\boldsymbol{1}_P^{T}(\boldsymbol{d} - \boldsymbol{y})\cdot f'(\boldsymbol{u})$ | $b = b + \alpha(d_p - y_p)f'(u_p)$ |

## GD / SGD For Logistic Regression Neuron (Sigmoid Activation)

| GD | SGD |
|---|---|
| $(\boldsymbol{X}, \boldsymbol{d})$ | $(\boldsymbol{x}_p, d_p)$ |
| $J = -\sum_{p=1}^{P} d_p \log\big(f(u_p)\big) + (1-d_p)\log\big(1-f(u_p)\big)$ | $J_p = -d_p \log\big(f(u_p)\big) - (1-d_p)\log\big(1-f(u_p)\big)$ |
| $\boldsymbol{u} = \boldsymbol{X}\boldsymbol{w} + b\boldsymbol{1}_P$ | $u_p = \boldsymbol{w}^{T}\boldsymbol{x}_p + b$ |
| $f(\boldsymbol{u}) = \dfrac{1}{1+e^{-\boldsymbol{u}}}$ | $f(u_p) = \dfrac{1}{1+e^{-u_p}}$ |
| $\boldsymbol{y} = 1(f(\boldsymbol{u}) > 0.5)$ | $y_p = 1(f(u_p) > 0.5)$ |
| $\boldsymbol{w} \leftarrow \boldsymbol{w} + \alpha\boldsymbol{X}^{T}(\boldsymbol{d} - f(\boldsymbol{u}))$ | $\boldsymbol{w} \leftarrow \boldsymbol{w} + \alpha\big(d_p - f(u_p)\big)\boldsymbol{x}_p$ |
| $b \leftarrow b + \alpha\boldsymbol{1}_P^{T}(\boldsymbol{d} - f(\boldsymbol{u}))$ | $b \leftarrow b + \alpha\big(d_p - f(u_p)\big)$ |

LAYERS

## SGD / GD For a Single layer of neurons

| Learning a layer of neurons | |
|---|---|
| **SGD** | $\boldsymbol{W} = \boldsymbol{W} - \alpha\boldsymbol{x}(\nabla_{\boldsymbol{u}}J)^{T}$  <br> $\boldsymbol{b} = \boldsymbol{b} - \alpha\nabla_{\boldsymbol{u}}J$ |
| **GD** | $\boldsymbol{W} = \boldsymbol{W} - \alpha\boldsymbol{X}^{T}\,\nabla_{\boldsymbol{U}}J$  <br> $\boldsymbol{b} = \boldsymbol{b} - \alpha(\nabla_{\boldsymbol{U}}J)^{T}\boldsymbol{1}_P$ |

## GD / SGD For a perceptron (sigmoid activation) layer - compare Perceptron!

| GD | SGD |
|---|---|
| $(\boldsymbol{X}, \boldsymbol{D})$ | $(\boldsymbol{x}, \boldsymbol{d})$ |
| $\boldsymbol{U} = \boldsymbol{X}\boldsymbol{W} + \boldsymbol{B}$ | $\boldsymbol{u} = \boldsymbol{W}^{T}\boldsymbol{x} + \boldsymbol{b}$ |
| $\boldsymbol{Y} = f(\boldsymbol{U})$ | $\boldsymbol{y} = f(\boldsymbol{u})$ |
| $\nabla_{\boldsymbol{U}}J = -(\boldsymbol{D} - \boldsymbol{Y})\cdot f'(\boldsymbol{U})$ | $\nabla_{\boldsymbol{u}}J = -(\boldsymbol{d} - \boldsymbol{y})\cdot f'(\boldsymbol{u})$ |
| $\boldsymbol{W} = \boldsymbol{W} - \alpha\boldsymbol{X}^{T}\,\nabla_{\boldsymbol{U}}J$ | $\boldsymbol{W} = \boldsymbol{W} - \alpha\boldsymbol{x}(\nabla_{\boldsymbol{u}}J)^{T}$ |
| $\boldsymbol{b} = \boldsymbol{b} - \alpha(\nabla_{\boldsymbol{U}}J)^{T}\boldsymbol{1}_P$ | $\boldsymbol{b} = \boldsymbol{b} - \alpha\nabla_{\boldsymbol{u}}J$ |

## GD / SGD For Softmax Layer - Compare Logistic Regression (discrete output at end)

| GD | SGD |
|---|---|
| $(\boldsymbol{X}, \boldsymbol{D})$ | $(\boldsymbol{x}, d)$ |
| $\boldsymbol{U} = \boldsymbol{X}\boldsymbol{W} + \boldsymbol{B}$ | $\boldsymbol{u} = \boldsymbol{W}^{T}\boldsymbol{x} + \boldsymbol{b}$ |
| $f(\boldsymbol{U}) = \dfrac{e^{\boldsymbol{U}}}{\sum_{k'=1}^{K} e^{\boldsymbol{U}_{k'}}}$ | $f(\boldsymbol{u}) = \dfrac{e^{u_k}}{\sum_{k'=1}^{K} e^{u_{k'}}}$ |
| $\boldsymbol{y} = \underset{k}{\mathrm{argmax}}\, f(\boldsymbol{U})$ | $y = \underset{k}{\mathrm{argmax}}\, f(\boldsymbol{u})$ |
| $\nabla_{\boldsymbol{U}}J = -(\boldsymbol{K} - f(\boldsymbol{U}))$ | $\nabla_{\boldsymbol{u}}J = -(1(\boldsymbol{k} = d) - f(\boldsymbol{u}))$ |
| $\boldsymbol{W} = \boldsymbol{W} - \alpha\boldsymbol{X}^{T}\,\nabla_{\boldsymbol{U}}J$ | $\boldsymbol{W} = \boldsymbol{W} - \alpha\boldsymbol{x}(\nabla_{\boldsymbol{u}}J)^{T}$ |
| $\boldsymbol{b} = \boldsymbol{b} - \alpha(\nabla_{\boldsymbol{U}}J)^{T}\boldsymbol{1}_P$ | $\boldsymbol{b} = \boldsymbol{b} - \alpha\nabla_{\boldsymbol{u}}J$ |

$\dfrac{d}{du}\,sigmoid(u) = y(1 - y).\ \ y = \dfrac{1}{1+e^{-u}}$

$\dfrac{d}{du}\,tanh(u) = 1 - y^2,\ \ y = \dfrac{e^{u}-e^{-u}}{e^{u}+e^{-u}}$

$$w_i \sim U(-a, a), \quad a = gain \times \sqrt{\frac{6}{n_{in}+n_{out}}}$$

$$w_i \sim N(0, \frac{gain}{n}), \quad n = num\ weights\ in\ neuron$$

Gain scale depending on the neuron activation function

| activation | linear | sigmoid | Tanh | ReLU | Leaky ReLU |
|---|---|---|---|---|---|
| gain | 1 | 1 | 5/3 | $\sqrt{2}$ | $\sqrt{\dfrac{1}{1+slope^2}}$ |

GD / SGD in Two Layer FFN (two layers, V weights and c bias for second layer)

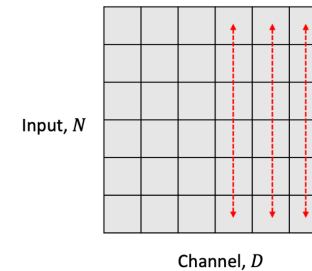| GD | SGD |
|---|---|
| $(X, D)$ | $(x, d)$ |
| $Z = XW + B$ | $z = W^T x + b$ |
| $H = g(Z)$ | $h = g(z)$ |
| $U = HV + C$ | $u = V^T h + c$ |
| $Y = f(U)$ | $y = f(u)$ |
| | |
| $\nabla_U J = \begin{cases} -(D - Y) \\ -(K - f(U)) \end{cases}$ | $\nabla_u J = \begin{cases} -(d - y) \\ (1(k = d) - f(u)) \end{cases}$ |
| $\nabla_Z J = (\nabla_U J) V^T \cdot g'(Z)$ | $\nabla_z J = V \nabla_u J \cdot g'(z)$ |
| | |
| $W \leftarrow W - \alpha X^T \nabla_Z J$ | $W \leftarrow W - \alpha x (\nabla_Z J)^T$ |
| $b \leftarrow b - \alpha (\nabla_Z J)^T \mathbf{1}_P$ | $b \leftarrow b - \alpha \nabla_z J$ |
| $V \leftarrow V - \alpha H^T \nabla_U J$ | $V \leftarrow V - \alpha h (\nabla_u J)^T$ |
| $c \leftarrow c - \alpha (\nabla_U J)^T \mathbf{1}_P$ | $c \leftarrow c - \alpha \nabla_u J$ |

Adam - Combines RMSprop and moment methods
Suggested defaults alpha=0.001, rho1=0.9, rho2=0.999, eps=10^-8

**Momentum term:**

$$s \leftarrow \rho_1 s + (1 - \rho_1) \nabla_W J$$

**Learning rate term:**

$$r \leftarrow \rho_2 r + (1 - \rho_2)(\nabla_W J)^2$$

$$s \leftarrow \frac{s}{1 - \rho_1}$$

$$r \leftarrow \frac{r}{1 - \rho_2}$$

$$W \leftarrow W - \frac{\alpha}{\varepsilon + \sqrt{r}} \cdot s$$

Batch Normalization



$$\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j}$$ 
Per-channel mean, shape is $1 \times D$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^{N} (x_{i,j} - \mu_j)^2$$ 
Per-channel variance, shape is $1 \times D$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$ 
Normalize the values, shape is $N \times D$
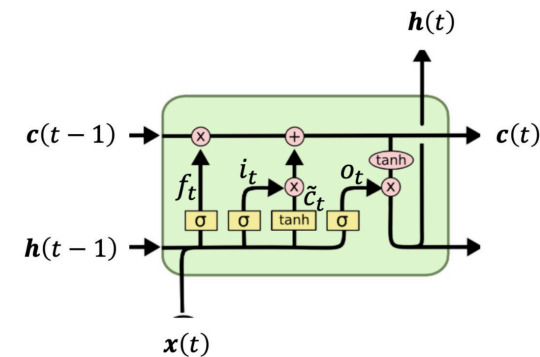
$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$ 
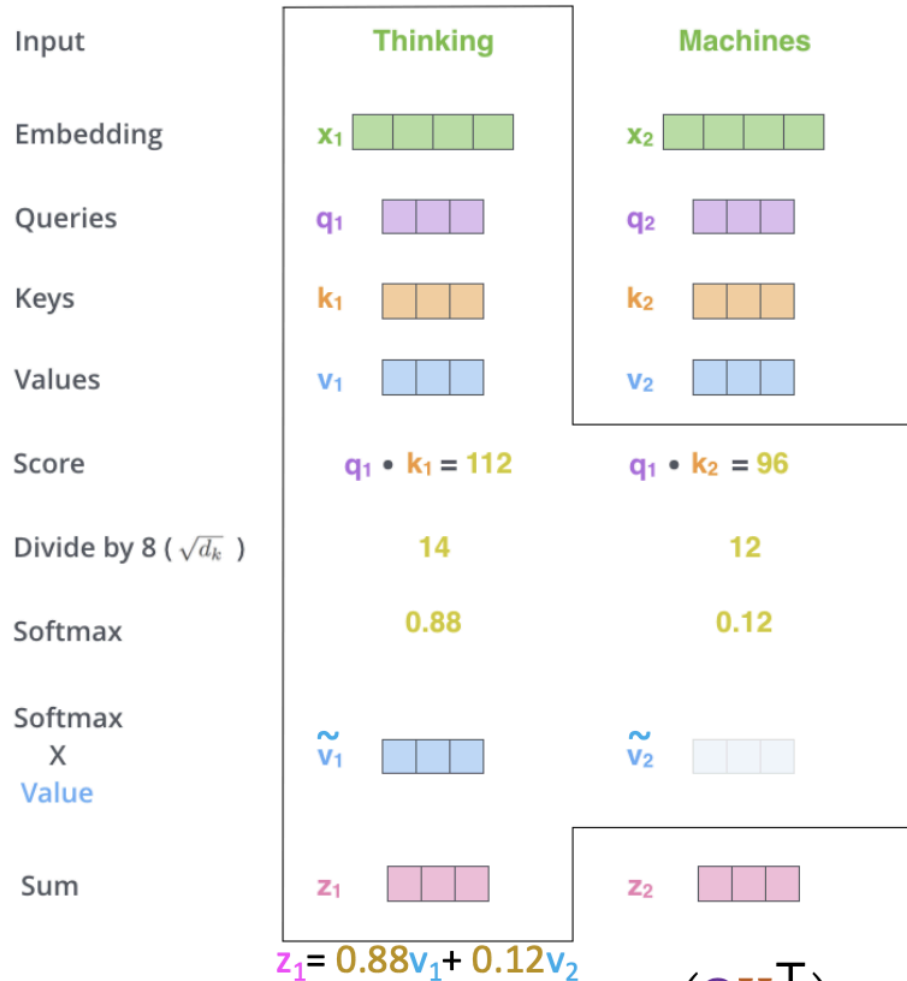**Scale and shift** the **normalized** values to add flexibility, output shape is $N \times D$

Learnable parameters: $\gamma$ and $\beta$, shape is $1 \times D$
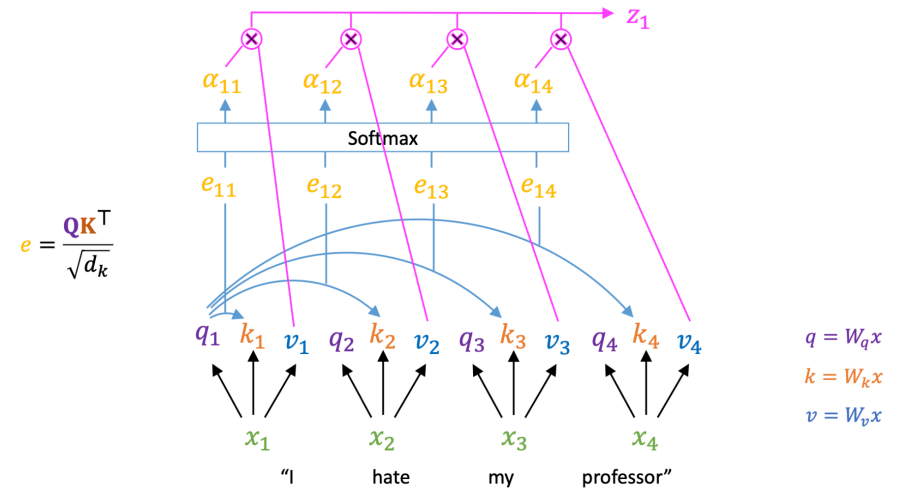
LSTM Cell

Key of LSTM is "state"
- A persistent module called the cell-state
- "State" is a representation of past history
- It comprises a common thread through time

| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \bullet k_1 = 112$ | $q_1 \bullet k_2 = 96$ |
| Divide by 8 ($\sqrt{d_k}$) | 14 | 12 |
| Softmax | 0.88 | 0.12 |
| Softmax X Value | $\tilde{v}_1$ | $\tilde{v}_2$ |
| Sum | $z_1$ | $z_2$ |

$$z_1 = 0.88v_1 + 0.12v_2$$

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{SoftMax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right)\mathbf{V}$$



$$e = \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}$$

$q = W_q x$
$k = W_k x$
$v = W_v x$

$x_1$ "I  $x_2$ hate  $x_3$ my  $x_4$ professor"

### First Step

Calculate the Query, Key, and Value matrices.

Pack our embeddings into a matrix **X**, and multiplying it by the weight matrices we've trained (**W$^Q$**, **W$^K$**, **W$^V$**)

Every row in the **X** matrix corresponds to a word in the input sentence.

### Second Step

Calculate the outputs of the self-attention layer. SoftMax is row-wise

$$\text{softmax}\left(\frac{Q \times K^\top}{\sqrt{d_k}}\right) V = Z$$