

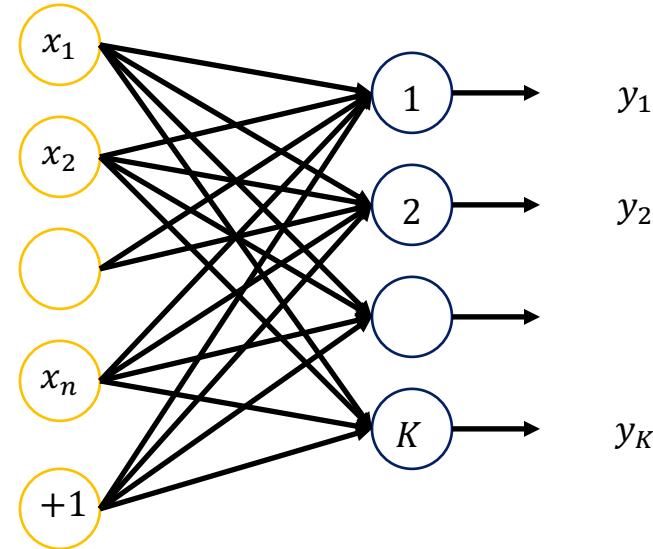
Chapter 4

Neuron Layers

Neural networks and deep learning

Weight matrix of layer

Consider a layer of K neurons.



Let \mathbf{w}_k and b_k denote the weight vector and bias of k th neuron.

Weights connected a neuron layer is represented by a weight matrix \mathbf{W} where columns are given by weight vectors connected to individual neurons:

$$\mathbf{W} = (\mathbf{w}_1 \quad \mathbf{w}_2 \quad \cdots \quad \mathbf{w}_K)$$

And a bias vector \mathbf{b} where each element corresponds to a bias of a neuron:

$$\mathbf{b} = (b_1, b_2, \dots, b_K)^T$$

Activation at layer for single input

Given an input pattern $\mathbf{x} \in \mathbb{R}^n$ to a layer of K neurons.

Synaptic input to k th neuron u_k :

$$u_k = \mathbf{w}_k^T \mathbf{x} + b_k$$

\mathbf{w}_k and b_k denote the weight vector and bias of k th neuron.

Synaptic input \mathbf{u} to the layer :

$$\mathbf{u} = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_K \end{pmatrix} = \begin{pmatrix} \mathbf{w}_1^T \mathbf{x} + b_1 \\ \mathbf{w}_2^T \mathbf{x} + b_2 \\ \vdots \\ \mathbf{w}_K^T \mathbf{x} + b_K \end{pmatrix} = \begin{pmatrix} \mathbf{w}_1^T \\ \mathbf{w}_2^T \\ \vdots \\ \mathbf{w}_K^T \end{pmatrix} \mathbf{x} + \mathbf{b} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$$

where \mathbf{W} is the weight matrix and \mathbf{b} is the bias vector of the layer.

Activation at layer for single input

Synaptic input to the layer:

$$\mathbf{u} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$$

The activation at the layer:

$$f(\mathbf{u}) = \begin{pmatrix} f(u_1) \\ f(u_2) \\ \vdots \\ f(u_K) \end{pmatrix}$$

Where $f(u_k)$ is the activation of k th neuron.

Synaptic input to a layer: batch input

Given a set $\{\mathbf{x}_p\}_{p=1}^P$ input patterns to a layer of K neurons where $\mathbf{x}_p \in \mathbb{R}^n$.

Synaptic input \mathbf{u}_p to the layer for an input pattern \mathbf{x}_p :

$$\mathbf{u}_p = \mathbf{W}^T \mathbf{x}_p + \mathbf{b}$$

The matrix \mathbf{U} of synaptic inputs to the layer for P patterns:

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$$

$$\mathbf{U} = \begin{pmatrix} \mathbf{u}_1^T \\ \mathbf{u}_2^T \\ \vdots \\ \mathbf{u}_P^T \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^T \mathbf{W} + \mathbf{b}^T \\ \mathbf{x}_2^T \mathbf{W} + \mathbf{b}^T \\ \vdots \\ \mathbf{x}_P^T \mathbf{W} + \mathbf{b}^T \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_P^T \end{pmatrix} \mathbf{W} + \begin{pmatrix} \mathbf{b}^T \\ \mathbf{b}^T \\ \vdots \\ \mathbf{b}^T \end{pmatrix} = \mathbf{XW} + \mathbf{B}$$

where rows of \mathbf{U} are synaptic inputs corresponding to individual input patterns.

The matrix $\mathbf{B} = \begin{pmatrix} \mathbf{b}^T \\ \mathbf{b}^T \\ \vdots \\ \mathbf{b}^T \end{pmatrix}$ has bias vector propagated as rows.

Synaptic input to a layer: batch input

The synaptic input to the layer due to a batch of patterns:

$$\mathbf{U} = \mathbf{X} \mathbf{W} + \mathbf{B}$$

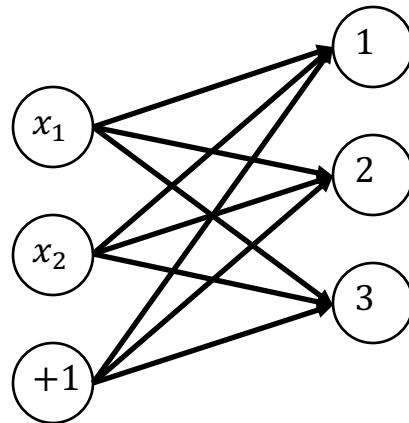
where rows of \mathbf{U} corresponds to synaptic inputs of the layer, corresponding to individual input patterns:

Activation of the layer:

$$f(\mathbf{U}) = \begin{pmatrix} f(\mathbf{u}_1^T) \\ f(\mathbf{u}_2^T) \\ \vdots \\ f(\mathbf{u}_P^T) \end{pmatrix} = \begin{pmatrix} f(\mathbf{u}_1)^T \\ f(\mathbf{u}_2)^T \\ \vdots \\ f(\mathbf{u}_P)^T \end{pmatrix}$$

where activation of each pattern is written as rows.

Example 1



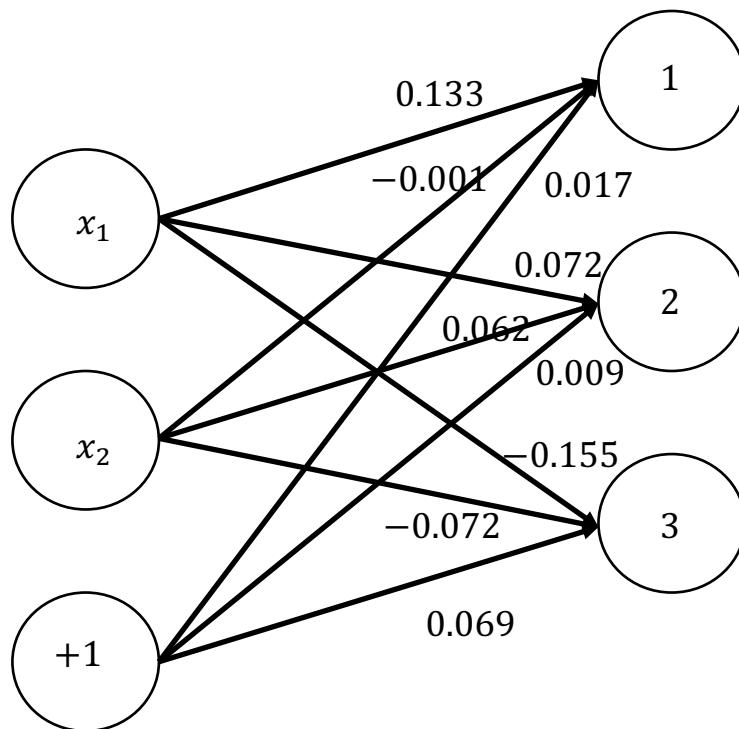
A perceptron layer of 3 neurons shown in the figure receives 2-dimensional inputs $(x_1, x_2)^T$, and has a weight matrix \mathbf{W} and a bias vector \mathbf{b} given by

$$\mathbf{W} = \begin{pmatrix} 0.133 & 0.072 & -0.155 \\ -0.001 & 0.062 & -0.072 \end{pmatrix} \text{ and } \mathbf{b} = \begin{pmatrix} 0.017 \\ 0.009 \\ 0.069 \end{pmatrix}$$

Using batch processing, find the output for input patterns:
 $\begin{pmatrix} 0.5 \\ -1.0 \\ -1.66 \end{pmatrix}$, $\begin{pmatrix} 0.78 \\ -0.51 \\ -0.65 \end{pmatrix}$, and $\begin{pmatrix} 0.04 \\ -0.2 \end{pmatrix}$.

Example 1

$$W = \begin{pmatrix} 0.133 & 0.072 & -0.155 \\ -0.001 & 0.062 & -0.072 \end{pmatrix} \text{ and } b = \begin{pmatrix} 0.017 \\ 0.009 \\ 0.069 \end{pmatrix}.$$



Example 1

$$\mathbf{W} = \begin{pmatrix} 0.133 & 0.072 & -0.155 \\ -0.001 & 0.062 & -0.072 \end{pmatrix} \text{ and } \mathbf{B} = \begin{pmatrix} 0.017 & 0.009 & 0.069 \\ 0.017 & 0.009 & 0.069 \\ 0.017 & 0.009 & 0.069 \\ 0.017 & 0.009 & 0.069 \end{pmatrix}.$$

Input as a batch of four patterns:

$$\mathbf{X} = \begin{pmatrix} 0.5 & -1.66 \\ -1.0 & -0.51 \\ 0.78 & -0.65 \\ 0.04 & -0.2 \end{pmatrix}$$

The synaptic input to the layer:

$$\begin{aligned} \mathbf{U} &= \mathbf{X}\mathbf{W} + \mathbf{B} \\ &= \begin{pmatrix} 0.5 & -1.66 \\ -1.0 & -0.51 \\ 0.78 & -0.65 \\ 0.04 & -0.2 \end{pmatrix} \begin{pmatrix} 0.133 & 0.072 & -0.155 \\ -0.001 & 0.062 & -0.072 \end{pmatrix} + \begin{pmatrix} 0.017 & 0.009 & 0.069 \\ 0.017 & 0.009 & 0.069 \\ 0.017 & 0.009 & 0.069 \\ 0.017 & 0.009 & 0.069 \end{pmatrix} \\ &= \begin{pmatrix} 0.085 & -0.059 & 0.111 \\ -0.115 & 0.094 & 0.26 \\ 0.121 & 0.024 & -0.005 \\ 0.022 & -0.001 & 0.077 \end{pmatrix} \end{aligned}$$

Example 1

$$\mathbf{U} = \begin{pmatrix} 0.085 & -0.059 & 0.111 \\ -0.115 & 0.094 & 0.26 \\ 0.121 & 0.024 & -0.005 \\ 0.022 & -0.001 & 0.077 \end{pmatrix}$$

For a perceptron layer

$$y = f(\mathbf{U}) = \frac{1}{1 + e^{-\mathbf{U}}} = \begin{pmatrix} 0.521 & 0.485 & 0.527 \\ 0.471 & 0.476 & 0.565 \\ 0.530 & 0.506 & 0.499 \\ 0.506 & 0.500 & 0.519 \end{pmatrix}$$

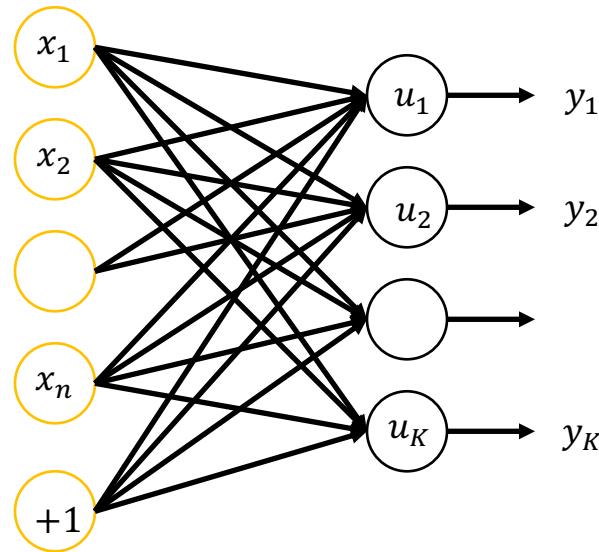
For example, third row corresponding to input

$$\mathbf{x} = \begin{pmatrix} 0.78 \\ -0.65 \end{pmatrix}$$

And output

$$\mathbf{y} = \begin{pmatrix} 0.530 \\ 0.506 \\ 0.499 \end{pmatrix}$$

Single layer of neurons

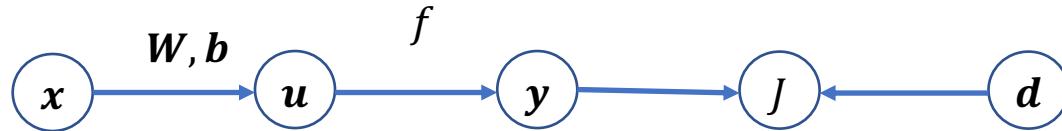


A single layer of K neurons processing an input of n -dimensions, connected by a weight matrix $\mathbf{W} = (\mathbf{w}_1 \quad \mathbf{w}_2 \quad \cdots \quad \mathbf{w}_K)$ and bias vector $\mathbf{b} = (b_1, b_2, \cdots b_K)^T$.

\mathbf{w}_k and b_k denote the weight vector and the bias of k th neuron.

SGD for single layer

Computational graph processing input (x, d) :



J denotes the cost function.

We need to compute gradients $\nabla_W J$ and $\nabla_b J$ to learn the weights W and biases b .

SGD for single layer

Consider k th neuron at the layer:

$$u_k = \mathbf{x}^T \mathbf{w}_k + b_k$$

And

$$\frac{\partial u_k}{\partial \mathbf{w}_k} = \mathbf{x}$$

The gradient of the cost with respect to the weight connected to k th neuron:

$$\nabla_{\mathbf{w}_k} J = \frac{\partial J}{\partial u_k} \frac{\partial u_k}{\partial \mathbf{w}_k} = \mathbf{x} \frac{\partial J}{\partial u_k} \quad (\text{A})$$

$$\nabla_{b_k} J = \frac{\partial J}{\partial u_k} \frac{\partial u_k}{\partial b_k} = \frac{\partial J}{\partial u_k} \quad (\text{B})$$

SGD for single layer

Gradient of J with respect to $\mathbf{W} = (\mathbf{w}_1 \quad \mathbf{w}_2 \quad \cdots \quad \mathbf{w}_K)$:

$$\begin{aligned}\nabla_{\mathbf{W}} J &= (\nabla_{\mathbf{w}_1} J \quad \nabla_{\mathbf{w}_2} J \quad \cdots \quad \nabla_{\mathbf{w}_K} J) \\ &= \left(\mathbf{x} \frac{\partial J}{\partial u_1} \quad \mathbf{x} \frac{\partial J}{\partial u_2} \quad \cdots \quad \mathbf{x} \frac{\partial J}{\partial u_K} \right) \quad \text{From (A)} \\ &= \mathbf{x} \left(\frac{\partial J}{\partial u_1} \quad \frac{\partial J}{\partial u_2} \quad \cdots \quad \frac{\partial J}{\partial u_K} \right) \\ &= \mathbf{x} (\nabla_{\mathbf{u}} J)^T \quad \text{(C)}\end{aligned}$$

where

$$\nabla_{\mathbf{u}} J = \begin{pmatrix} \frac{\partial J}{\partial u_1} \\ \frac{\partial J}{\partial u_2} \\ \vdots \\ \frac{\partial J}{\partial u_K} \end{pmatrix}$$

SGD for single layer

$$\nabla_{\mathbf{W}} J = \mathbf{x} (\nabla_{\mathbf{u}} J)^T$$

Similarly from (B):

$$\nabla_{\mathbf{b}} J = \begin{pmatrix} \frac{\partial J}{\partial b_1} \\ \frac{\partial J}{\partial b_2} \\ \vdots \\ \frac{\partial J}{\partial b_K} \end{pmatrix} = \begin{pmatrix} \frac{\partial J}{\partial u_1} \\ \frac{\partial J}{\partial u_2} \\ \vdots \\ \frac{\partial J}{\partial u_K} \end{pmatrix} = \nabla_{\mathbf{u}} J \quad (\text{D})$$

SGD for single layer

That is, from (C) and (D), by computing gradient $\nabla_w J$ with respect to synaptic input, the gradient of cost J with respect to the weights and biases can be obtained:

$$\begin{aligned}\nabla_w J &= x(\nabla_u J)^T \\ \nabla_b J &= \nabla_u J\end{aligned}$$

SGD for single layer

Substituting (C) and (D) in gradient descent equations:

$$\begin{aligned}\mathbf{W} &\leftarrow \mathbf{W} - \alpha \nabla_{\mathbf{W}} J \\ \mathbf{b} &\leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{b}} J\end{aligned}$$

SGD equations for single layer:

$$\begin{aligned}\mathbf{W} &\leftarrow \mathbf{W} - \alpha x (\nabla_{\mathbf{w}} J)^T \\ \mathbf{b} &\leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{w}} J\end{aligned} \tag{E}$$

GD for single layer

Given a set of patterns $\{(\mathbf{x}_p, \mathbf{d}_p)\}_{p=1}^P$ where $\mathbf{x}_p \in \mathbf{R}^n$ and $\mathbf{d}_p \in \mathbf{R}^K$ for regression and $d_p \in \{1, 2, \dots, K\}$ for classification.

The cost J is given by the sum of cost due to individual patterns:

$$J = \sum_{p=1}^P J_p$$

Where Then,

$$\nabla_{\mathbf{W}} J = \sum_{p=1}^P \nabla_{\mathbf{W}} J_p$$

GD for single layer

Substituting $\nabla_{\mathbf{W}} J_p = \mathbf{x}_p (\nabla_{\mathbf{u}_p} J_p)^T$ from (C) :

$$\begin{aligned}
 \nabla_{\mathbf{W}} J &= \sum_{p=1}^P \mathbf{x}_p (\nabla_{\mathbf{u}_p} J_p)^T \\
 &= \sum_{p=1}^P \mathbf{x}_p (\nabla_{\mathbf{u}_p} J)^T && \text{since } \nabla_{\mathbf{u}_p} J = \nabla_{\mathbf{u}_p} J_p. \\
 &= \mathbf{x}_1 (\nabla_{\mathbf{u}_1} J)^T + \mathbf{x}_2 (\nabla_{\mathbf{u}_2} J)^T + \dots + \mathbf{x}_P (\nabla_{\mathbf{u}_P} J)^T \\
 &= (\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_P) \begin{pmatrix} (\nabla_{\mathbf{u}_1} J)^T \\ (\nabla_{\mathbf{u}_2} J)^T \\ \vdots \\ (\nabla_{\mathbf{u}_P} J)^T \end{pmatrix} \\
 &= \mathbf{X}^T \nabla_{\mathbf{U}} J
 \end{aligned} \tag{F}$$

Note that $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_P^T \end{pmatrix}$ and $\mathbf{U} = \begin{pmatrix} \mathbf{u}_1^T \\ \mathbf{u}_2^T \\ \vdots \\ \mathbf{u}_P^T \end{pmatrix}$

GD for single layer

$$J = \sum_{p=1}^P J_p$$

$$\text{Then, } \nabla_{\mathbf{u}_p} J = \nabla_{\mathbf{u}_p} J_p. \quad (\text{G})$$

$$\begin{aligned}\nabla_{\mathbf{b}} J &= \sum_{p=1}^P \nabla_{\mathbf{b}} J_p \\ &= \sum_{p=1}^P \nabla_{\mathbf{u}_p} J_p && \text{Substituting from (D)} \\ &= \sum_{p=1}^P \nabla_{\mathbf{u}_p} J && \text{Substituting from (G)} \\ &= \nabla_{\mathbf{u}_1} J + \nabla_{\mathbf{u}_2} J + \cdots + \nabla_{\mathbf{u}_P} J \\ &= (\nabla_{\mathbf{u}_1} J \quad \nabla_{\mathbf{u}_2} J \quad \cdots \quad \nabla_{\mathbf{u}_P} J) \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \\ &= (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P && (\text{H})\end{aligned}$$

where $\mathbf{1}_P = (1, 1, \dots, 1)^T$ is a vector of P ones.

GD for single layer

From (F) and (H):

$$\begin{aligned}\nabla_{\mathbf{W}} J &= \mathbf{X}^T \nabla_{\mathbf{U}} J \\ \nabla_{\mathbf{b}} J &= (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P\end{aligned}$$

That is, by computing gradient $\nabla_{\mathbf{U}} J$ with respect to synaptic input, the weights and biases can be updated.

Substituting in gradient descent equations:

$$\begin{aligned}\mathbf{W} &\leftarrow \mathbf{W} - \alpha \nabla_{\mathbf{W}} J \\ \mathbf{b} &\leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{b}} J\end{aligned}$$

GD equations for single layer networks:

$$\begin{aligned}\mathbf{W} &\leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{U}} J \\ \mathbf{b} &\leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P\end{aligned}\tag{I}$$

Learning a single layer

Gradient Descent Learning

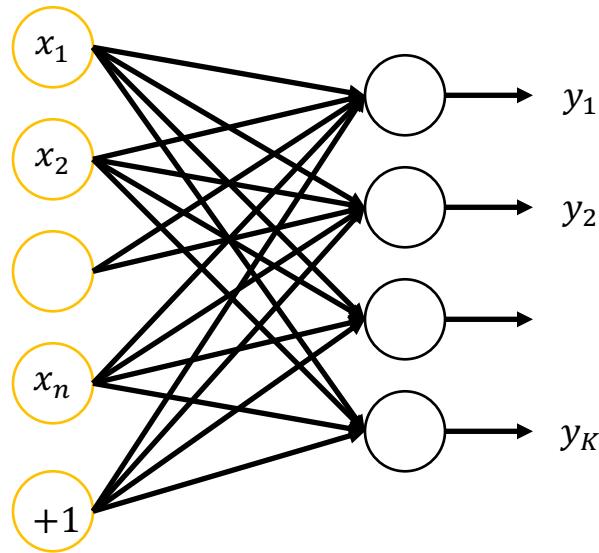
$$\begin{aligned} \mathbf{W} &\leftarrow \mathbf{W} - \alpha \nabla_{\mathbf{W}} J \\ \mathbf{b} &\leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{b}} J \end{aligned}$$

Learning of a layer

SGD	$\nabla_{\mathbf{W}} J = \mathbf{x}(\nabla_{\mathbf{u}} J)^T$ $\nabla_{\mathbf{b}} J = \nabla_{\mathbf{u}} J$	$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{x}(\nabla_{\mathbf{u}} J)^T$ $\mathbf{b} \leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{u}} J$
GD	$\nabla_{\mathbf{W}} J = \mathbf{X}^T \nabla_{\mathbf{U}} J$ $\nabla_{\mathbf{b}} J = (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P$	$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{U}} J$ $\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P$

For learning, what is necessary is to compute $\nabla_{\mathbf{u}} J$ for SGD and to compute $\nabla_{\mathbf{U}} J$ for GD. These gradients are dependent on the type of the layer.

Perceptron layer



A layer of perceptrons performs **multidimensional regression**.

A layer of K perceptrons learns a multidimensional non-linear mapping:
 $\phi: \mathbf{R}^n \rightarrow \mathbf{R}^K$

SGD for perceptron layer

Given a training pattern (\mathbf{x}, \mathbf{d})

where $\mathbf{x} = (x_1, x_2, \dots, x_K) \in \mathbf{R}^n$ and $\mathbf{d} = (d_1, d_2, \dots, d_K)^T \in \mathbf{R}^K$.

The square-error cost function:

$$J = \frac{1}{2} \sum_{k=1}^K (d_k - y_k)^2$$

where $y_k = f(u_k)$ and $u_k = \mathbf{x}^T \mathbf{w}_k + b_k$.

Gradient of J with respect to u_k :

$$\frac{\partial J}{\partial u_k} = \frac{\partial J}{\partial y_k} \frac{\partial y_k}{\partial u_k} = -(d_k - y_k) \frac{\partial y_k}{\partial u_k} = -(d_k - y_k) f'(u_k)$$

SGD for perceptron layer

Substituting in gradient J with respect to $\mathbf{u} = (u_1, u_2, \dots, u_K)^T$:

$$\nabla_{\mathbf{u}} J = \begin{pmatrix} \nabla_{u_1} J \\ \nabla_{u_2} J \\ \vdots \\ \nabla_{u_K} J \end{pmatrix} = - \begin{pmatrix} (d_1 - y_1)f'(u_1) \\ (d_2 - y_2)f'(u_2) \\ \vdots \\ (d_K - y_K)f'(u_K) \end{pmatrix} = -(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})$$

$$\text{where } \mathbf{d} = \begin{pmatrix} d_1 \\ d_2 \\ \vdots \\ d_K \end{pmatrix}, \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_K \end{pmatrix}, f'(\mathbf{u}) = \begin{pmatrix} f'(u_1) \\ f'(u_2) \\ \vdots \\ f'(u_K) \end{pmatrix}$$

and ‘ \cdot ’ denotes element-wise multiplication.

For a perceptron layer:

$$\nabla_{\mathbf{u}} J = -(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u}) \quad (\text{J})$$

SGD for perceptron layer

From (E), SGD learning for single layer:

$$\begin{aligned}\mathbf{W} &\leftarrow \mathbf{W} - \alpha \mathbf{x} (\nabla_{\mathbf{w}} J)^T \\ \mathbf{b} &\leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{b}} J\end{aligned}$$

Substituting form (J), SGD equations for perceptron:

$$\begin{aligned}\mathbf{W} &\leftarrow \mathbf{W} + \alpha \mathbf{x} ((\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u}))^T \\ \mathbf{b} &\leftarrow \mathbf{b} + \alpha (\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})\end{aligned}$$

SGD for perceptron layer

Given a training dataset $\{(x, d)\}$

Set learning parameter α

Initialize W and b

Repeat until convergence:

For every pattern (x, d) :

$$u = W^T x + b$$

$$y = f(u)$$

$$\nabla_u J = -(d - y) \cdot f'(u)$$

$$W \leftarrow W - \alpha x (\nabla_u J)^T$$

$$b \leftarrow b - \alpha \nabla_u J$$

GD for perceptron layer

Given a training dataset $\{(\mathbf{x}_p, \mathbf{d}_p)\}_{p=1}^P$

where $\mathbf{x}_p = (x_{p1}, x_{p2}, \dots, x_{pn})^T \in \mathbf{R}^n$ and $\mathbf{d}_p = (d_{p1}, d_{p2}, \dots, d_{pK})^T \in \mathbf{R}^K$.

The cost function J is given by the sum of square errors (s.s.e.):

$$J = \frac{1}{2} \sum_{p=1}^P \sum_{k=1}^K (d_{pk} - y_{pk})^2$$

J can be written as the sum of cost due to individual patterns:

$$J = \sum_{p=1}^P J_p$$

where $J_p = \frac{1}{2} \sum_{k=1}^K (d_{pk} - y_{pk})^2$ is the square error for the p th pattern.

GD for perceptron layer

$$\mathbf{U} = \begin{pmatrix} \mathbf{u}_1^T \\ \mathbf{u}_2^T \\ \vdots \\ \mathbf{u}_P^T \end{pmatrix} \rightarrow \nabla_{\mathbf{U}} J = \begin{pmatrix} (\nabla_{\mathbf{u}_1} J)^T \\ (\nabla_{\mathbf{u}_2} J)^T \\ \vdots \\ (\nabla_{\mathbf{u}_P} J)^T \end{pmatrix} = \begin{pmatrix} (\nabla_{\mathbf{u}_1} J_1)^T \\ (\nabla_{\mathbf{u}_2} J_2)^T \\ \vdots \\ (\nabla_{\mathbf{u}_P} J_P)^T \end{pmatrix}$$

From (J), substituting $\nabla_{\mathbf{u}} J = -(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})$:

$$\begin{aligned} \nabla_{\mathbf{U}} J &= - \begin{pmatrix} ((\mathbf{d}_1 - \mathbf{y}_1) \cdot f'(\mathbf{u}_1))^T \\ ((\mathbf{d}_2 - \mathbf{y}_2) \cdot f'(\mathbf{u}_2))^T \\ \vdots \\ ((\mathbf{d}_P - \mathbf{y}_P) \cdot f'(\mathbf{u}_P))^T \end{pmatrix} = - \begin{pmatrix} (\mathbf{d}_1^T - \mathbf{y}_1^T) \cdot f'(\mathbf{u}_1^T) \\ (\mathbf{d}_2^T - \mathbf{y}_2^T) \cdot f'(\mathbf{u}_2^T) \\ \vdots \\ (\mathbf{d}_P^T - \mathbf{y}_P^T) \cdot f'(\mathbf{u}_P^T) \end{pmatrix} \\ &= -(\mathbf{D} - \mathbf{Y}) \cdot f'(\mathbf{U}) \end{aligned}$$

where $\mathbf{D} = \begin{pmatrix} \mathbf{d}_1^T \\ \mathbf{d}_2^T \\ \vdots \\ \mathbf{d}_P^T \end{pmatrix}$, $\mathbf{Y} = \begin{pmatrix} \mathbf{y}_1^T \\ \mathbf{y}_2^T \\ \vdots \\ \mathbf{y}_P^T \end{pmatrix}$, and $\mathbf{U} = \begin{pmatrix} \mathbf{u}_1^T \\ \mathbf{u}_2^T \\ \vdots \\ \mathbf{u}_P^T \end{pmatrix}$

GD for perceptron layer

For a perceptron layer (batch input);

$$\nabla_{\mathbf{U}} J = -(\mathbf{D} - \mathbf{Y}) \cdot f'(\mathbf{U})$$

GD equation for a single layer from (I);

$$\begin{aligned}\mathbf{W} &\leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{U}} J \\ \mathbf{b} &\leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P\end{aligned}$$

Substituting $\nabla_{\mathbf{U}} J$, we get GD equations for a perceptron layer:

$$\begin{aligned}\mathbf{W} &\leftarrow \mathbf{W} + \alpha \mathbf{X}^T (\mathbf{D} - \mathbf{Y}) \cdot f'(\mathbf{U}) \\ \mathbf{b} &\leftarrow \mathbf{b} + \alpha ((\mathbf{D} - \mathbf{Y}) \cdot f'(\mathbf{U}))^T \mathbf{1}_P\end{aligned}$$

GD for perceptron layer

Given a training dataset (\mathbf{X}, \mathbf{D})

Set learning parameter α

Initialize \mathbf{W} and \mathbf{b}

Repeat until convergence:

$$\mathbf{U} = \mathbf{X}\mathbf{W} + \mathbf{B}$$

$$\mathbf{Y} = f(\mathbf{U}) = \frac{1}{1+e^{-\mathbf{U}}}$$

$$\nabla_{\mathbf{U}} J = -(\mathbf{D} - \mathbf{Y}) \cdot f'(\mathbf{U})$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{U}} J$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P$$

Learning a perceptron layer

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

Example 2

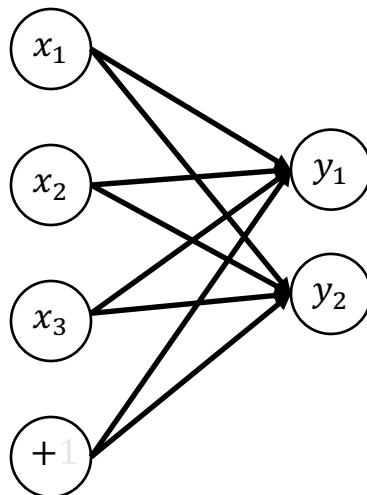
Design a perceptron layer to perform the following mapping using GD learning:

$x = (x_1, x_2, x_3)$	$y = (y_1, y_2)$
(0.77, 0.02, 0.63)	(0.37, 0.47)
(0.75, 0.50, 0.22)	(0.36, 0.38)
(0.20, 0.76, 0.17)	(0.35, 0.25)
(0.09, 0.69, 0.95)	(0.48, 0.42)
(0.00, 0.51, 0.81)	(0.36, 0.29)
(0.61, 0.72, 0.29)	(0.44, 0.52)
(0.92, 0.71, 0.54)	(0.60, 0.52)
(0.14, 0.37, 0.67)	(0.28, 0.37)

Use $\alpha = 0.1$.

Example 2

$$X = \begin{pmatrix} 0.77 & 0.02 & 0.63 \\ 0.75 & 0.50 & 0.22 \\ 0.20 & 0.76 & 0.17 \\ 0.09 & 0.69 & 0.95 \\ 0.00 & 0.51 & 0.81 \\ 0.61 & 0.72 & 0.29 \\ 0.92 & 0.71 & 0.54 \\ 0.14 & 0.37 & 0.67 \end{pmatrix} \text{ and } D = \begin{pmatrix} 0.37 & 0.47 \\ 0.36 & 0.38 \\ 0.35 & 0.25 \\ 0.48 & 0.42 \\ 0.36 & 0.29 \\ 0.44 & 0.52 \\ 0.60 & 0.52 \\ 0.28 & 0.37 \end{pmatrix}$$



Output $y_1, y_2 \in [0, 1]$

So, activation function for both neurons:

$$f(u) = \frac{1}{1 + e^{-u}}$$

$$f'(u) = y(1 - y)$$

Learning factor $\alpha = 0.1$.

Weights and biases are initialized:

$$W = \begin{pmatrix} 0.03 & 0.04 \\ 0.01 & 0.04 \\ 0.02 & 0.04 \end{pmatrix} \text{ and } b = \begin{pmatrix} 0.0 \\ 0.0 \end{pmatrix}$$

Example 2

Iteration 1:

$$\mathbf{U} = \mathbf{XW} + \mathbf{B}$$

$$\mathbf{U} = \begin{pmatrix} 0.77 & 0.02 & 0.63 \\ 0.75 & 0.50 & 0.22 \\ 0.20 & 0.76 & 0.17 \\ 0.09 & 0.69 & 0.95 \\ 0.00 & 0.51 & 0.81 \\ 0.61 & 0.72 & 0.29 \\ 0.92 & 0.71 & 0.54 \\ 0.14 & 0.37 & 0.67 \end{pmatrix} \begin{pmatrix} 0.03 & 0.04 \\ 0.01 & 0.04 \\ 0.02 & 0.04 \end{pmatrix} + \begin{pmatrix} 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \end{pmatrix} = \begin{pmatrix} 0.03 & 0.06 \\ 0.03 & 0.06 \\ 0.02 & 0.05 \\ 0.03 & 0.07 \\ 0.02 & 0.05 \\ 0.03 & 0.07 \\ 0.04 & 0.09 \\ 0.02 & 0.05 \end{pmatrix}$$

$$Y = f(\mathbf{U}) = \frac{1}{1+e^{-\mathbf{U}}} = \begin{pmatrix} 0.51 & 0.51 \\ 0.51 & 0.52 \\ 0.50 & 0.51 \\ 0.51 & 0.52 \\ 0.50 & 0.51 \\ 0.51 & 0.52 \\ 0.51 & 0.52 \\ 0.50 & 0.51 \end{pmatrix}$$

$$\text{Mean square error} = \frac{1}{8} \sum_{p=1}^8 \sum_{k=1}^2 (d_{pk} - y_{pk})^2 = \frac{1}{8} \sum_{p=1}^8 (d_{p1} - y_{p1})^2 + (d_{p2} - y_{p2})^2 = 0.04$$

Example 2

$$f'(\mathbf{U}) = \mathbf{Y} \cdot (1 - \mathbf{Y}) = \begin{pmatrix} 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \end{pmatrix}$$

$$\nabla_{\mathbf{U}} J = -(\mathbf{D} - \mathbf{Y}) \cdot f'(\mathbf{U})$$

$$= - \left(\begin{pmatrix} 0.37 & 0.47 \\ 0.36 & 0.38 \\ 0.35 & 0.25 \\ 0.48 & 0.42 \\ 0.36 & 0.29 \\ 0.44 & 0.52 \\ 0.60 & 0.52 \\ 0.28 & 0.37 \end{pmatrix} - \begin{pmatrix} 0.51 & 0.51 \\ 0.51 & 0.52 \\ 0.50 & 0.51 \\ 0.51 & 0.52 \\ 0.50 & 0.51 \\ 0.51 & 0.52 \\ 0.51 & 0.52 \\ 0.50 & 0.51 \end{pmatrix} \right) \cdot \begin{pmatrix} 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \\ 0.25 & 0.25 \end{pmatrix} = \begin{pmatrix} 0.03 & 0.01 \\ 0.04 & 0.03 \\ 0.04 & 0.07 \\ 0.01 & 0.03 \\ 0.04 & 0.06 \\ 0.02 & 0.00 \\ -0.02 & 0.00 \\ 0.06 & 0.04 \end{pmatrix}$$

Example 2

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{U}} J$$

$$\mathbf{W} = \begin{pmatrix} 0.03 & 0.04 \\ 0.01 & 0.04 \\ 0.02 & 0.04 \end{pmatrix} - 0.1 \begin{pmatrix} 0.77 & 0.02 & 0.63 \\ 0.75 & 0.50 & 0.22 \\ 0.20 & 0.76 & 0.17 \\ 0.09 & 0.69 & 0.95 \\ 0.00 & 0.51 & 0.81 \\ 0.61 & 0.72 & 0.29 \\ 0.92 & 0.71 & 0.54 \\ 0.14 & 0.37 & 0.67 \end{pmatrix}^T \begin{pmatrix} 0.03 & 0.01 \\ 0.04 & 0.03 \\ 0.04 & 0.07 \\ 0.01 & 0.03 \\ 0.04 & 0.06 \\ 0.02 & 0.00 \\ -0.02 & 0.00 \\ 0.06 & 0.04 \end{pmatrix} = \begin{pmatrix} 0.03 & 0.04 \\ 0.01 & 0.04 \\ 0.02 & 0.04 \end{pmatrix}$$

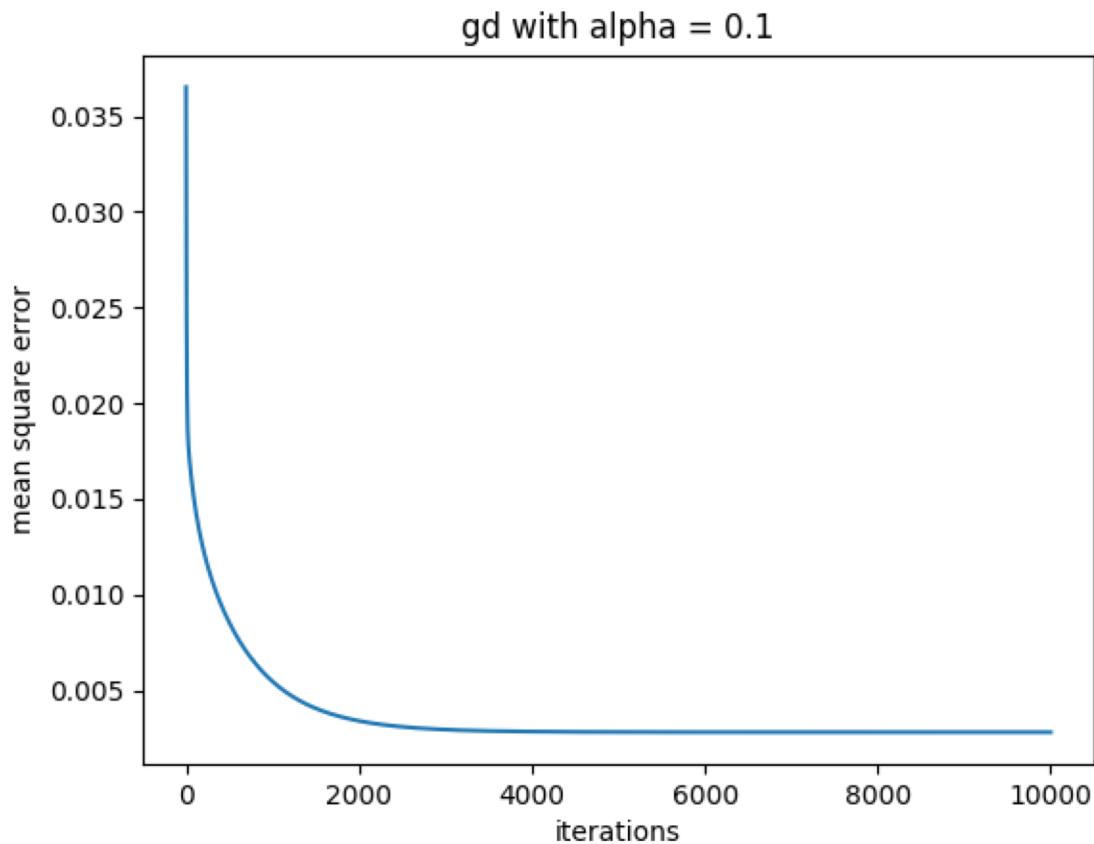
$$\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P$$

$$\mathbf{b} = \begin{pmatrix} 0.0 \\ 0.0 \end{pmatrix} - 0.1 \begin{pmatrix} 0.03 & 0.01 \\ 0.03 & 0.03 \\ 0.04 & 0.06 \\ 0.00 & 0.02 \\ 0.03 & 0.05 \\ 0.01 & 0.00 \\ -0.02 & 0.00 \\ 0.05 & 0.03 \end{pmatrix}^T \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 0.0 \\ 0.0 \end{pmatrix}$$

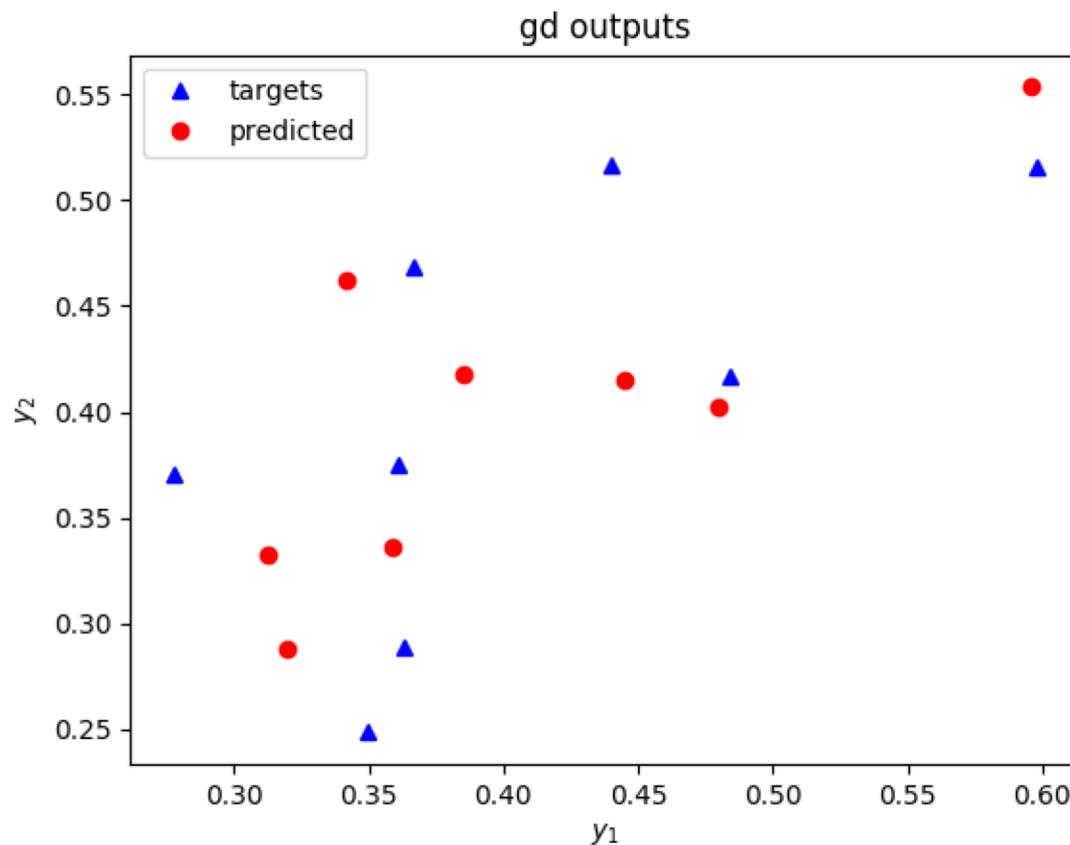
Example 2

iter	Y	mse	W	b
2	$\begin{pmatrix} 0.50 & 0.51 \\ 0.50 & 0.51 \\ 0.50 & 0.50 \\ 0.50 & 0.51 \\ 0.50 & 0.50 \\ 0.50 & 0.51 \\ 0.50 & 0.51 \\ 0.50 & 0.50 \end{pmatrix}$	0.04	$\begin{pmatrix} 0.02 & 0.04 \\ 0.00 & 0.03 \\ 0.01 & 0.03 \end{pmatrix}$	$\begin{pmatrix} -0.02 \\ -0.02 \end{pmatrix}$
10000	$\begin{pmatrix} 0.34 & 0.46 \\ 0.39 & 0.42 \\ 0.32 & 0.29 \\ 0.48 & 0.40 \\ 0.36 & 0.34 \\ 0.45 & 0.42 \\ 0.59 & 0.55 \\ 0.31 & 0.33 \end{pmatrix}$	0.003	$\begin{pmatrix} 1.06 & 1.12 \\ 1.39 & 0.38 \\ 1.11 & 0.82 \end{pmatrix}$	$\begin{pmatrix} -2.2 \\ -1.54 \end{pmatrix}$

Example 2

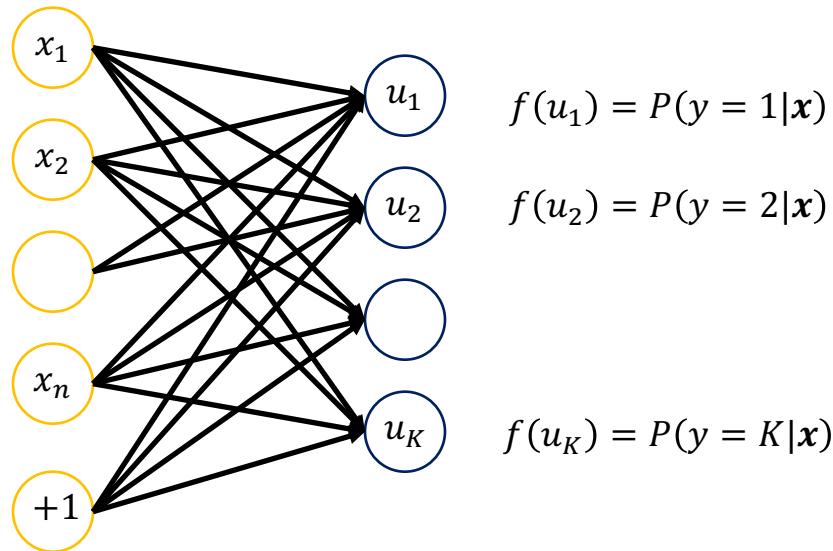


Example 2



Softmax layer

Softmax layer is the extension of logistic regression to **multiclass classification** problem, which is also known as *multinomial logistic regression*.



Every neuron in the softmax layer corresponds to one class label. The activation of the neuron gives the probability of the input belonging to class label.

Softmax layer

The K neurons in the softmax layer performs multinomial logistic regression and represent K classes.

The activation of each neuron k estimates the probability $P(y = k|x)$ that the input \mathbf{x} belongs to the class k :

$$P(y = k|x) = f(u_k) = \frac{e^{u_k}}{\sum_{k'=1}^K e^{u_{k'}}}$$

where $u_k = \mathbf{x}^T \mathbf{w}_k + b_k$, and \mathbf{w}_k is weight vector and b_k is bias of neuron k .

The above activation function f is known as **softmax activation function**.

Softmax layer

The output y denotes the class label of the input pattern, which is given by

$$y = \operatorname{argmax}_k P(y = k | \mathbf{x}) = \operatorname{argmax}_k f(u_k)$$

That is, the class label is assigned to the class with the maximum activation.

SGD for softmax layer

Given a training pattern (\mathbf{x}, d) where $\mathbf{x} \in \mathbf{R}^n$ and $d \in \{1, 2, \dots, K\}$.

The SGD cost function for learning is by the *multiclass cross-entropy*:

$$J = - \sum_{k=1}^K 1(d = k) \log(f(u_k))$$

where $\mathbf{u} = (u_1, u_2, \dots, u_K)^T$ is synaptic input to the layer and u_k is the synaptic input to the k the neuron.

The cost function can also be written as

$$J = -\log(f(u_d))$$

where d is the target label of input \mathbf{x} .

SGD for softmax layer

$$J = -\log(f(u_d))$$

The gradient with respect to u_k is given by

$$\frac{\partial J}{\partial u_k} = -\frac{1}{f(u_d)} \frac{\partial f(u_d)}{\partial u_k} \quad (\text{K})$$

where

$$\frac{\partial f(u_d)}{\partial u_k} = \frac{\partial}{\partial u_k} \left(\frac{e^{u_d}}{\sum_{k'=1}^K e^{u_{k'}}} \right)$$

SGD for softmax layer

If $k = d$:

$$\begin{aligned}\frac{\partial f(u_d)}{\partial u_k} &= \frac{\partial}{\partial u_k} \left(\frac{e^{u_k}}{\sum_{k'=1}^K e^{u_{k'}}} \right) \\ &= \frac{(\sum_{k'=1}^K e^{u_{k'}}) e^{u_k} - e^{u_k} e^{u_k}}{(\sum_{k'=1}^K e^{u_{k'}})^2} \\ &= \frac{e^{u_k}}{\sum_{k'=1}^K e^{u_{k'}}} \left(1 - \frac{e^{u_k}}{\sum_{k'=1}^K e^{u_{k'}}} \right) \\ &= f(u_k)(1 - f(u_k)) \\ &= f(u_d)(1 - f(u_k))\end{aligned}\tag{L}$$

If $k \neq d$:

$$\begin{aligned}\frac{\partial f(u_d)}{\partial u_k} &= \frac{\partial}{\partial u_k} \left(\frac{e^{u_d}}{\sum_{k'=1}^K e^{u_{k'}}} \right) && 1(d=k) \\ &= -\frac{e^{u_d} e^{u_k}}{(\sum_{k'=1}^K e^{u_{k'}})^2} \\ &= -f(u_d)f(u_k)\end{aligned}\tag{M}$$

SGD for softmax layer

Combining (L) and (M):

$$\frac{\partial f(u_d)}{\partial u_k} = f(u_d)(1(d = k) - f(u_k))$$

Substituting in (K):

$$\frac{\partial J}{\partial u_k} = -(1(d = k) - f(u_k))$$

Gradient J with respect to $\mathbf{u} = (u_1, u_2, \dots, u_K)^T$:

$$\nabla_{\mathbf{u}} J = \begin{pmatrix} \nabla_{u_1} J \\ \nabla_{u_2} J \\ \vdots \\ \nabla_{u_K} J \end{pmatrix} = - \begin{pmatrix} 1(d = 1) - f(u_1) \\ 1(d = 2) - f(u_2) \\ \vdots \\ 1(d = K) - f(u_K) \end{pmatrix} = -(1(\mathbf{k} = d) - f(\mathbf{u}))$$

where $\mathbf{k} = (1 \quad 2 \quad \dots \quad K)^T$

SGD for softmax layer

For a softmax layer:

$$\nabla_{\mathbf{u}} J = -(1(\mathbf{k} = d) - f(\mathbf{u}))$$

Where $1(\mathbf{k} = d) = \begin{pmatrix} 1(d = 1) \\ 1(d = 2) \\ \vdots \\ 1(d = K) \end{pmatrix}$ and $f(\mathbf{u}) = \begin{pmatrix} f(u_1) \\ f(u_2) \\ \vdots \\ f(u_K) \end{pmatrix}$

SGD learning for single layer :

$$\begin{aligned}\mathbf{W} &\leftarrow \mathbf{W} - \alpha \mathbf{x} (\nabla_{\mathbf{u}} J)^T \\ \mathbf{b} &\leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{u}} J\end{aligned}$$

Substituting $\nabla_{\mathbf{u}} J$, we get SGD equations for softmax layer:

$$\begin{aligned}\mathbf{W} &\leftarrow \mathbf{W} + \alpha \mathbf{x} (1(\mathbf{k} = d) - f(\mathbf{u}))^T \\ \mathbf{b} &\leftarrow \mathbf{b} + \alpha (1(\mathbf{k} = d) - f(\mathbf{u}))\end{aligned}$$

SGD for softmax layer

Given a training dataset $\{(x, d)\}$

Set learning parameter α

Initialize W and b

Repeat until convergence:

For every pattern (x, d) :

$$u = W^T x + b$$

$$f(u) = \frac{e^{u_k}}{\sum_{k'=1}^K e^{u_{k'}}}$$

$$\nabla_u J = -(1(k=d) - f(u))$$

$$W \leftarrow W - \alpha x (\nabla_u J)^T$$

$$b \leftarrow b - \alpha \nabla_w J$$

GD for softmax layer

Given a set of patterns $\{(x_p, d_p)\}_{p=1}^P$ where $x_p \in \mathbb{R}^n$ and $d_p \in \{1, 2, \dots, K\}$.

The cost function of the *softmax layer* is given by the *multiclass cross-entropy*:

$$J = - \sum_{p=1}^P \left(\sum_{k=1}^K 1(d_p = k) \log(f(u_{pk})) \right)$$

where u_{pk} is the synaptic input to the k the neuron for input x_p .

The cost function J can also be written as

$$J = - \sum_{p=1}^P \log(f(u_{pd_p}))$$

where d_p is the target of input x_p .

GD for softmax layer

J can be written as the sum of cost due to individual patterns:

$$J = \sum_{p=1}^P J_p$$

where $J_p = -\log(f(u_{pd_p}))$ is the cross-entropy for the p th pattern.

GD for softmax layer

$$\begin{aligned}\nabla_{\mathbf{U}} J &= \begin{pmatrix} (\nabla_{\mathbf{u}_1} J_1)^T \\ (\nabla_{\mathbf{u}_2} J_2)^T \\ \vdots \\ (\nabla_{\mathbf{u}_P} J_P)^T \end{pmatrix} \\ &= - \begin{pmatrix} (1(\mathbf{k} = d_1) - f(\mathbf{u}_1))^T \\ (1(\mathbf{k} = d_2) - f(\mathbf{u}_2))^T \\ \vdots \\ (1(\mathbf{k} = d_K) - f(\mathbf{u}_K))^T \end{pmatrix} \quad \text{Substituting from (N)} \\ &= -(\mathbf{K} - f(\mathbf{U}))\end{aligned}$$

$$\text{where } \mathbf{K} = \begin{pmatrix} 1(\mathbf{k} = d_1)^T \\ 1(\mathbf{k} = d_2)^T \\ \vdots \\ 1(\mathbf{k} = d_P)^T \end{pmatrix}.$$

GD for softmax layer

For a softmax layer (batch input):

$$\nabla_{\mathbf{U}} J = -(\mathbf{K} - f(\mathbf{U}))$$

Substituting $\nabla_{\mathbf{U}} J$ in gradient descent equations:

$$\begin{aligned}\mathbf{W} &\leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{U}} J \\ \mathbf{b} &\leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P\end{aligned}$$

we get GD equations for a softmax layer:

$$\begin{aligned}\mathbf{W} &\leftarrow \mathbf{W} + \alpha \mathbf{X}^T (\mathbf{K} - f(\mathbf{U})) \\ \mathbf{b} &\leftarrow \mathbf{b} + \alpha (\mathbf{K} - f(\mathbf{U}))^T \mathbf{1}_P\end{aligned}$$

GD for softmax layer

Given training set (X, d)

Set learning rate α

Initialize W and b

Iterate until convergence:

$$U = XW + B$$

$$f(U) = \frac{e^U}{\sum_{k=1}^K e^{U_k}}$$

$$\nabla_U J = -(K - f(U))$$

$$W \leftarrow W - \alpha X^T \nabla_U J$$

$$b \leftarrow b - \alpha (\nabla_U J)^T \mathbf{1}_P$$

Learning a softmax layer

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

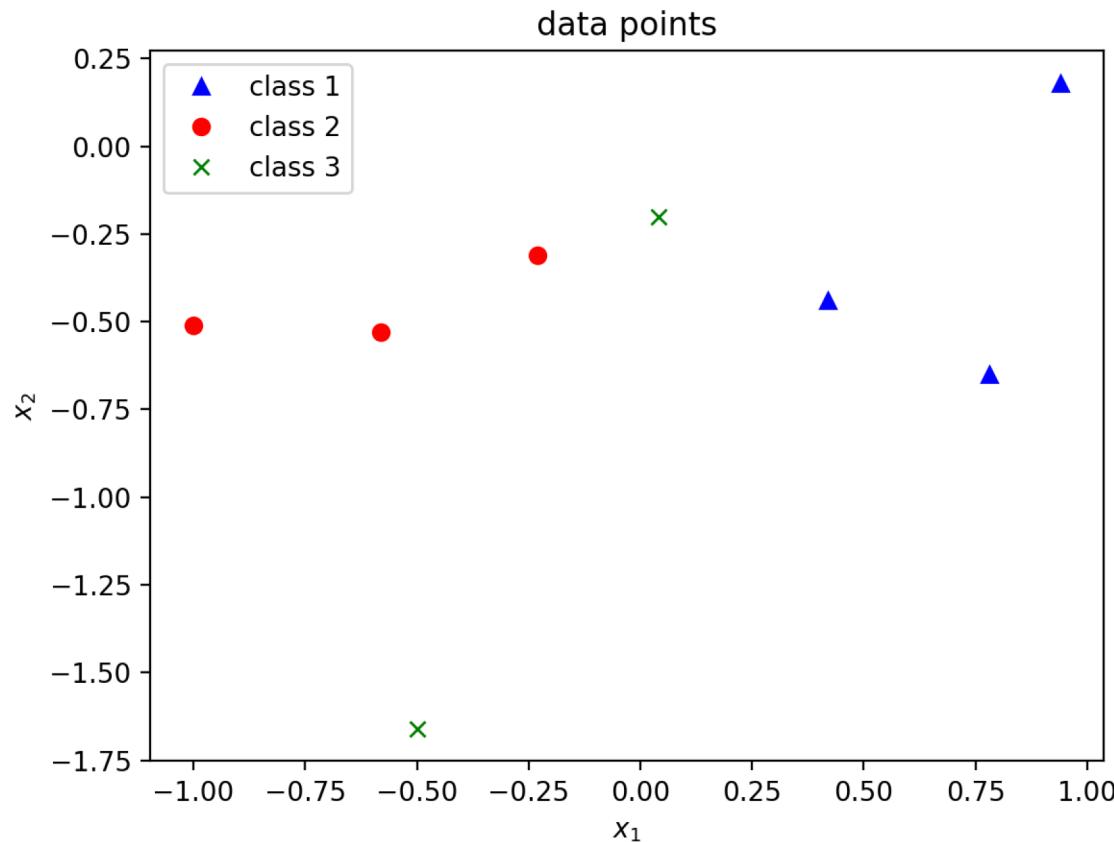
Example 3

Train a softmax regression layer of neurons to perform the following classification:

- (0.94 0.18) → *class A*
- (−0.58 −0.53) → *class B*
- (−0.23 −0.31) → *class B*
- (0.42 −0.44) → *class A*
- (0.5 −1.66) → *class C*
- (−1.0 −0.51) → *class B*
- (0.78 −0.65) → *class A*
- (0.04 −0.20) → *class C*

Use a learning factor $\alpha = 0.05$.

Example 3

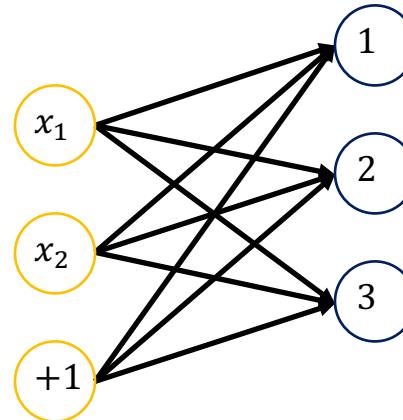


Example 3

Let $y = \begin{cases} 1, & \text{for class } A \\ 2, & \text{for class } B \\ 3, & \text{for class } C \end{cases}$

$$X = \begin{pmatrix} 0.94 & 0.18 \\ -0.58 & -0.53 \\ -0.23 & -0.31 \\ 0.42 & -0.44 \\ 0.5 & -1.66 \\ -1.0 & -0.51 \\ 0.78 & -0.65 \\ 0.04 & -0.2 \end{pmatrix}, \quad \boldsymbol{d} = \begin{pmatrix} 1 \\ 2 \\ 2 \\ 1 \\ 3 \\ 2 \\ 1 \\ 3 \end{pmatrix}$$

$$\boldsymbol{K} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$



Example 3

Initialize $\mathbf{W} = \begin{pmatrix} 0.77 & 0.02 & 0.63 \\ 0.75 & 0.50 & 0.23 \end{pmatrix}$, $\mathbf{b} = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix}$,

Example 3

$$U = \begin{pmatrix} 0.86 & 0.11 & 0.64 \\ -0.84 & -0.28 & -0.49 \\ -0.41 & -0.16 & -0.22 \\ -0.01 & -0.21 & 0.17 \\ -0.86 & -0.82 & -0.06 \\ -1.15 & -0.27 & -0.75 \\ 0.11 & -0.31 & 0.35 \\ -0.12 & -0.10 & -0.02 \end{pmatrix}$$

$$f(u_{12}) = \frac{e^{0.11}}{e^{0.86} + e^{0.11} + e^{0.64}}$$

$$f(U) = \frac{e^{(U)}}{\sum_{k=1}^K e^{(U)}} = \begin{pmatrix} 0.44 & 0.21 & 0.35 \\ 0.24 & 0.42 & 0.34 \\ 0.29 & 0.37 & 0.35 \\ 0.33 & 0.27 & 0.40 \\ 0.23 & 0.24 & 0.52 \\ 0.20 & 0.49 & 0.31 \\ 0.34 & 0.22 & 0.43 \\ 0.32 & 0.33 & 0.35 \end{pmatrix}$$

Example 3

$$\mathbf{y} = \underset{k}{\operatorname{argmax}} \{f(\mathbf{U})\} = \begin{pmatrix} 1 \\ 2 \\ 2 \\ 3 \\ 3 \\ 2 \\ 3 \\ 3 \end{pmatrix}$$

$$\text{Errors} = - \sum_{p=1}^8 1(\mathbf{d} \neq \mathbf{y}) = 2$$

$$\begin{aligned} \text{Entropy, } J(\mathbf{W}, \mathbf{b}) &= - \sum_{p=1}^8 \log \left(f(u_{pd_p}) \right) \\ &= -\log(0.44) - \log(0.42) - \cdots - \log(0.35) \\ &= 7.26 \end{aligned}$$

Example 3

$$\nabla_{\mathbf{U}} J = -(\mathbf{K} - f(\mathbf{U})) = - \left(\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} - \begin{pmatrix} 0.44 & 0.21 & 0.35 \\ 0.24 & 0.42 & 0.34 \\ 0.29 & 0.37 & 0.35 \\ 0.33 & 0.27 & 0.40 \\ 0.23 & 0.24 & 0.52 \\ 0.20 & 0.49 & 0.31 \\ 0.34 & 0.22 & 0.43 \\ 0.32 & 0.33 & 0.35 \end{pmatrix} \right) = \begin{pmatrix} -0.56 & 0.21 & 0.35 \\ 0.24 & -0.58 & 0.34 \\ 0.29 & -0.63 & 0.35 \\ -0.67 & 0.27 & 0.40 \\ 0.23 & 0.24 & -0.48 \\ 0.20 & -0.51 & 0.31 \\ -0.65 & 0.22 & 0.43 \\ 0.32 & 0.33 & -0.65 \end{pmatrix}$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{U}} J$$

$$\mathbf{W} = \begin{pmatrix} 0.77 & 0.02 & 0.63 \\ 0.75 & 0.50 & 0.23 \end{pmatrix} - 0.05 \begin{pmatrix} 0.94 & 0.18 \\ -0.58 & -0.53 \\ -0.23 & -0.31 \\ 0.42 & -0.44 \\ 0.5 & -1.66 \\ -1.0 & -0.51 \\ 0.78 & -0.65 \\ 0.04 & -0.2 \end{pmatrix}^T \begin{pmatrix} -0.56 & 0.21 & 0.35 \\ 0.24 & -0.58 & 0.34 \\ 0.29 & -0.63 & 0.35 \\ -0.67 & 0.27 & 0.40 \\ 0.23 & 0.24 & -0.48 \\ 0.20 & -0.51 & 0.31 \\ -0.65 & 0.22 & 0.43 \\ 0.32 & 0.33 & -0.65 \end{pmatrix}$$

$$= \begin{pmatrix} 0.85 & -0.06 & 0.63 \\ 0.76 & 0.50 & 0.22 \end{pmatrix}$$

Example 3

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P$$

$$\mathbf{b} = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix} - 0.05 \begin{pmatrix} -0.56 & 0.21 & 0.35 \\ 0.24 & -0.58 & 0.34 \\ 0.29 & -0.63 & 0.35 \\ -0.67 & 0.27 & 0.40 \\ 0.23 & 0.24 & -0.48 \\ 0.20 & -0.51 & 0.31 \\ -0.65 & 0.22 & 0.43 \\ 0.32 & 0.33 & -0.65 \end{pmatrix}^T \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 0.03 \\ 0.02 \\ -0.05 \end{pmatrix}$$

Example 3

$$\mathbf{U} = \mathbf{XW} + \mathbf{B}$$

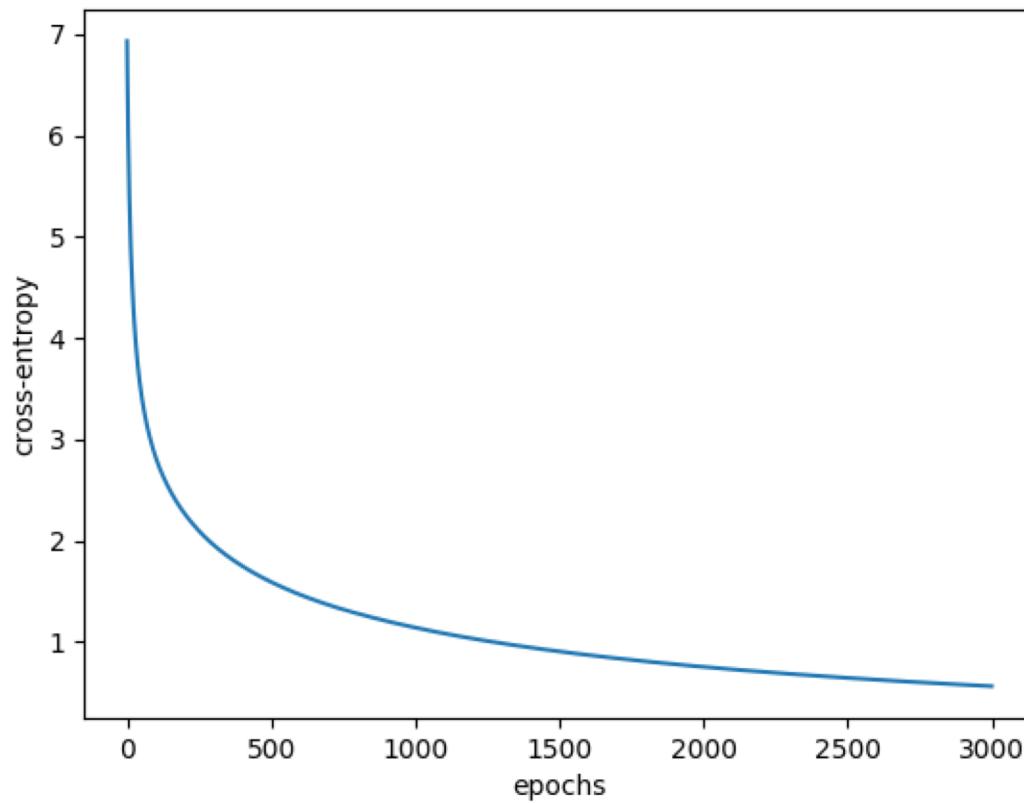
$$f(\mathbf{U}) = \frac{e^{\mathbf{U}}}{\sum_{k=1}^K e^{\mathbf{U}_k}}$$

$$\nabla_{\mathbf{U}} J = -\left(K - f(\mathbf{U})\right)$$

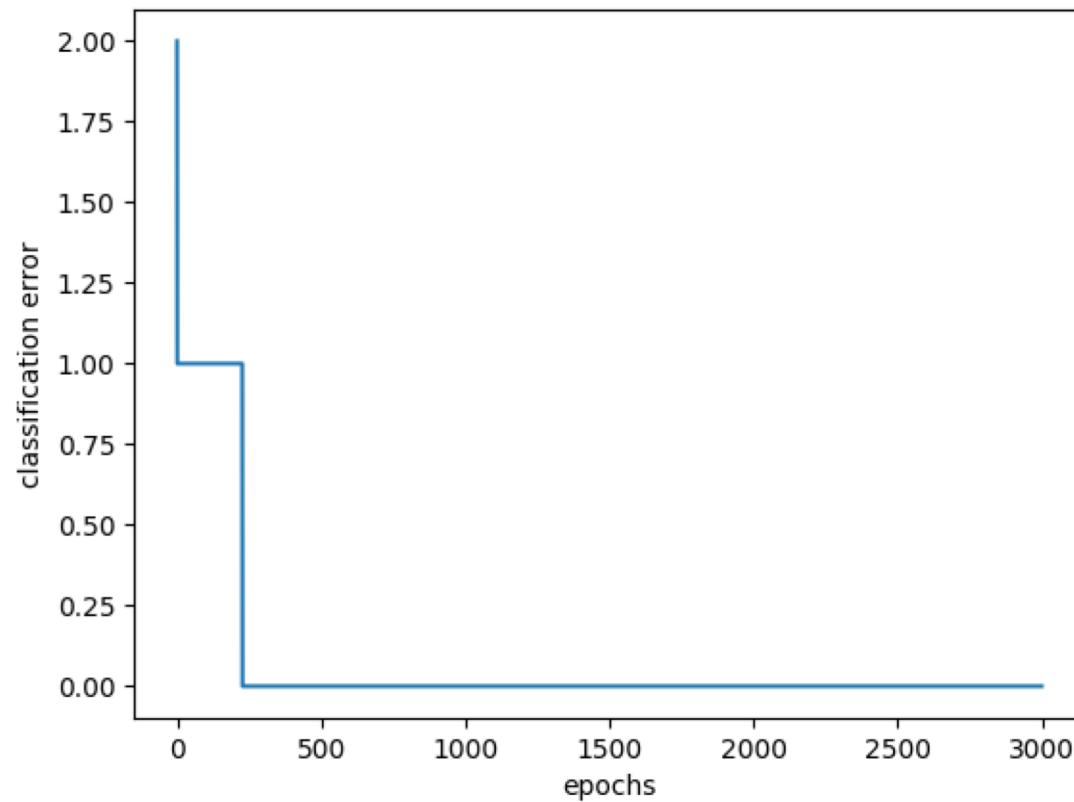
$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{U}} J$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P$$

Example 3



Example 3



Iris dataset

Iris dataset:

<https://archive.ics.uci.edu/ml/datasets/Iris>

Three classes of iris flower:



Setosa



Versicolour



Virginica

Four features:

Sepal length, sepal width, petal length, petal width

Iris dataset

150 data points, 50 for each class

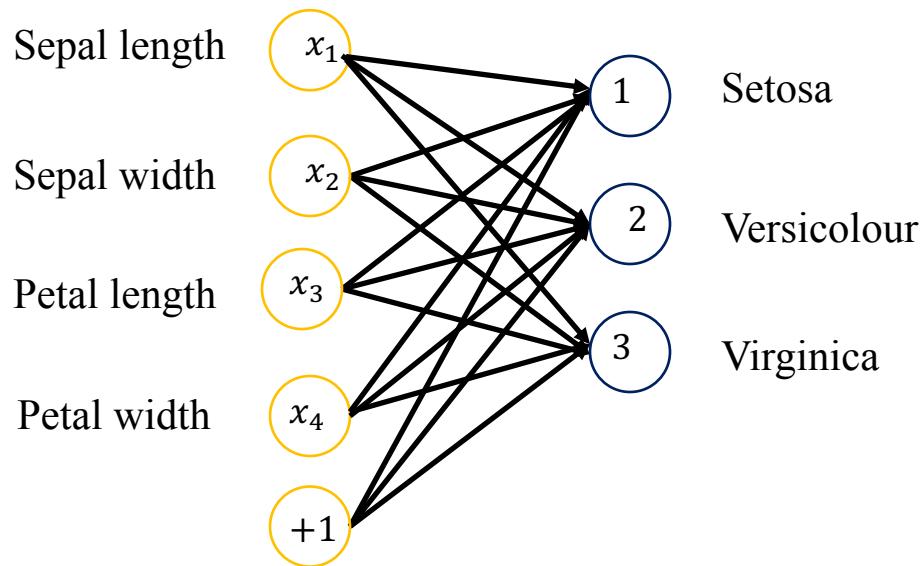
Features:

```
[ -7.4333333e-01  4.4600000e-01 -2.35866667e+00 -9.98666667e-01]
[ -9.4333333e-01 -5.4000000e-02 -2.35866667e+00 -9.98666667e-01]
[ -1.14333333e+00  1.4600000e-01 -2.45866667e+00 -9.98666667e-01]
```

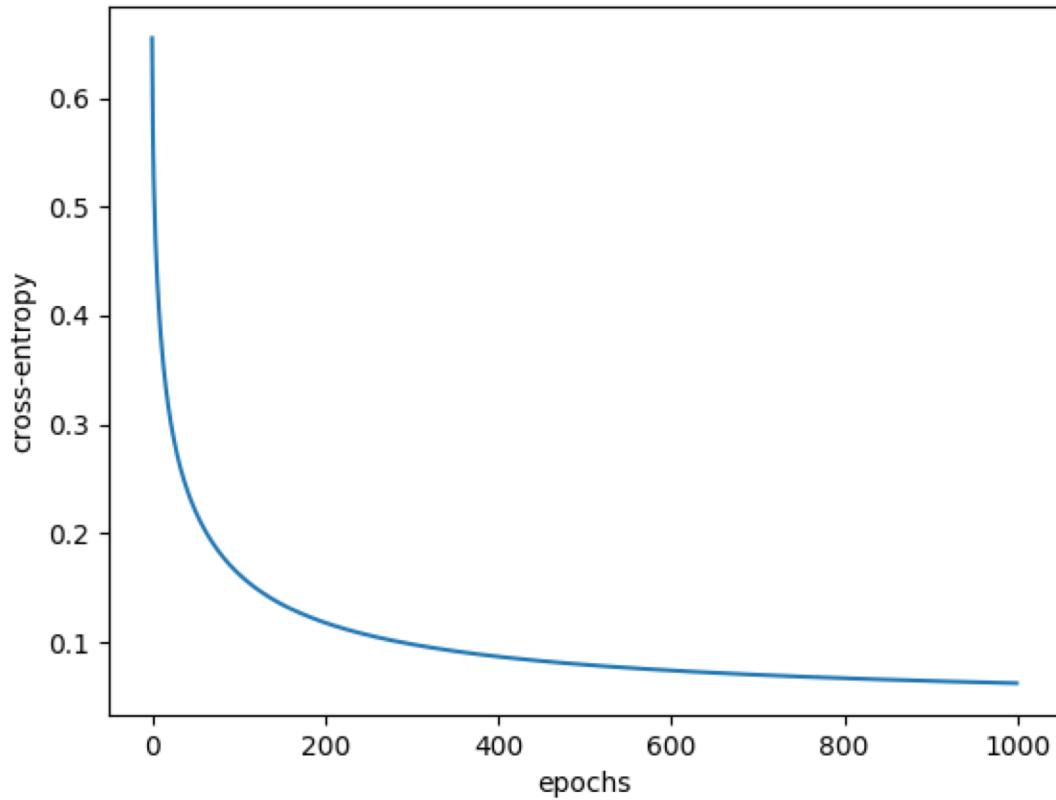
Labels:

```
[0 0 0 0 0 0 0 0 ..1 1 1 1 1 1 1 1 1 1 .... 2 2 2 2 2 2 2 2]
```

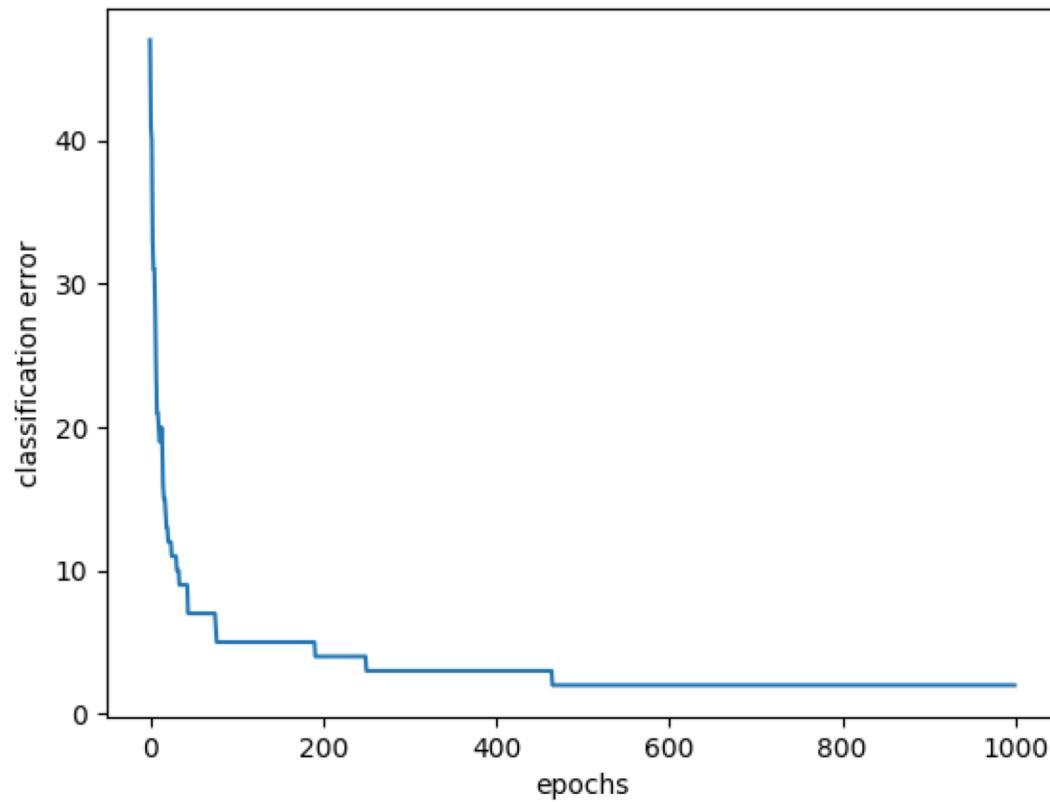
Example 4: Softmax classification of iris data



Example 4



Example 4



Final classification error = 2

Initialization of weights

Random initialization is inefficient

At initialization, it is desirable that weights

- are small and near zero to operate in the linear region of the activation function
- Preserve the variance of activation and gradients.

Initialization from a uniform distribution

For logistic sigmoid:

$$w \sim Uniform \left[-\frac{4\sqrt{6}}{\sqrt{n_{in} + n_{out}}}, \frac{4\sqrt{6}}{\sqrt{n_{in} + n_{out}}} \right]$$

For others:

$$w \sim Uniform \left[-\frac{\sqrt{6}}{\sqrt{n_{in} + n_{out}}}, \frac{\sqrt{6}}{\sqrt{n_{in} + n_{out}}} \right]$$

n_{in} is the input number of nodes and n_{out} is the number of neurons in the layer. *Uniform* is a uniformly distributed number within limits.

Initialization from a truncated normal distribution

$$w \sim \text{truncated_normal} \left[\text{mean} = 0, \text{std} = \frac{1}{\sqrt{n_{in}}} \right]$$

In the truncated normal, the samples that are two s.d. away from the center are discarded and resampled again.

Revision: neurons and layers

Classification	Perceptron	Logistic neurons
Two-class	Discrete perceptron	Logistic regression neuron
Multiclass	Discrete perceptron layer	Softmax neuron layer

Regression	Linear	Non-linear
One dimensional	Linear neuron	Perceptron
Multi-dimensional	Linear neuron layer	Perceptron layer

Summary: GD for layers

$$\begin{aligned}
 & (\mathbf{X}, \mathbf{D}) \\
 & \mathbf{U} = \mathbf{XW} + \mathbf{B} \\
 & \mathbf{W} = \mathbf{W} - \alpha \mathbf{X}^T (\nabla_{\mathbf{U}} J) \\
 & \mathbf{b} = \mathbf{b} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P
 \end{aligned}$$

layer	$f(\mathbf{U}), \mathbf{Y}$	$\nabla_{\mathbf{U}} J$
Linear neuron layer	$\mathbf{Y} = f(\mathbf{U}) = \mathbf{U}$	$-(\mathbf{D} - \mathbf{Y})$
Perceptron layer	$\mathbf{Y} = f(\mathbf{U}) = \frac{1}{1 + e^{-\mathbf{U}}}$	$-(\mathbf{D} - \mathbf{Y}) \cdot f'(\mathbf{U})$
Softmax layer	$f(\mathbf{U}) = \frac{e^{\mathbf{U}}}{\sum_{k=1}^K e^{\mathbf{U}_k}}$ $\mathbf{y} = \underset{k}{\operatorname{argmax}} f(\mathbf{U})$	$-(\mathbf{K} - f(\mathbf{U}))$

Summary: SGD for layers

$$\begin{aligned}
 & (\mathbf{x}, \mathbf{d}) \\
 & \mathbf{u} = \mathbf{W}^T \mathbf{x} + \mathbf{b} \\
 & \mathbf{W} = \mathbf{W} - \alpha \mathbf{x} (\nabla_{\mathbf{u}} J)^T \\
 & \mathbf{b} = \mathbf{b} - \alpha (\nabla_{\mathbf{u}} J)
 \end{aligned}$$

layer	$f(\mathbf{u}), \mathbf{y}$	$\nabla_{\mathbf{u}} J$
Linear neuron layer	$\mathbf{y} = f(\mathbf{u}) = \mathbf{u}$	$-(\mathbf{d} - \mathbf{y})$
Perceptron layer	$\mathbf{y} = f(\mathbf{u}) = \frac{1}{1 + e^{-\mathbf{u}}}$	$-(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})$
Softmax layer	$f(\mathbf{u}) = \frac{e^{\mathbf{u}}}{\sum_{k'=1}^K e^{k'}}$ $\mathbf{y} = \underset{k}{\operatorname{argmax}} f(\mathbf{u})$	$-(1(\mathbf{k} = d) - f(\mathbf{u}))$