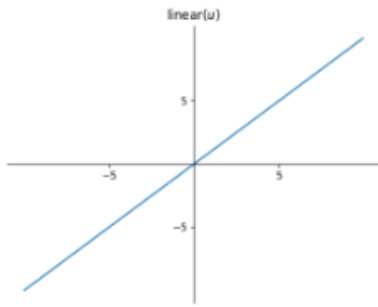


## Linear Neuron



A linear neuron is given by

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

The activation function is

$$y = f(u) = u$$

The output can be expanded to

$$y = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b, \text{ with } x_1 \dots \text{ as inputs}$$

The linear neuron performs linear regression and the weights and bias act as regression coefficients

Given a training dataset  $\{(x_p, d_p)\}_{p=1}^P$  where input  $\mathbf{x} \in \mathbf{R}^n$  and target  $d_p \in \mathbf{R}$  training a linear neuron finds a regression function  $\phi: \mathbf{R}^n \rightarrow \mathbf{R}$ , that is given by the linear mapping:  $y = \mathbf{w}^T \mathbf{x} + b$

Cost function:  $J = \frac{1}{2} (d - y)^2$

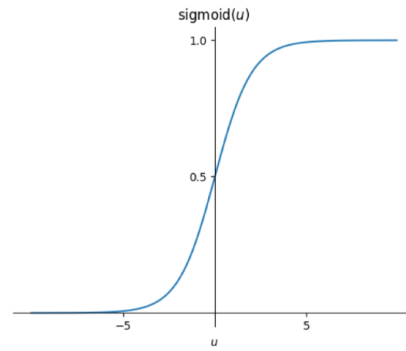
$$\nabla_u J = \frac{\partial J}{\partial u} = \frac{\partial J}{\partial y} = -(d - y)$$

$$\nabla_{\mathbf{w}} J = \frac{\partial J}{\partial \mathbf{w}} = \frac{\partial J}{\partial u} \frac{\partial u}{\partial \mathbf{w}} = -(d - y) \cdot \mathbf{x}$$

$$\nabla_b J = \frac{\partial J}{\partial b} = \frac{\partial J}{\partial u} \frac{\partial u}{\partial b} = -(d - y) \cdot 1$$

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

## Perceptron- logistic function



A Perceptron neuron is given by

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

The activation function is

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

The square error is used as cost function for learning

$\phi$  is a non-linear function and does non-linear regression

For the training date set  $\{(x_p, d_p)\}_{p=1}^P$  the continuous perceptron finds a functional mapping  $\phi: \mathbf{R}^n \rightarrow \mathbf{R}$  by learning from training data.

Cost function:  $J = \frac{1}{2} (d - y)^2$

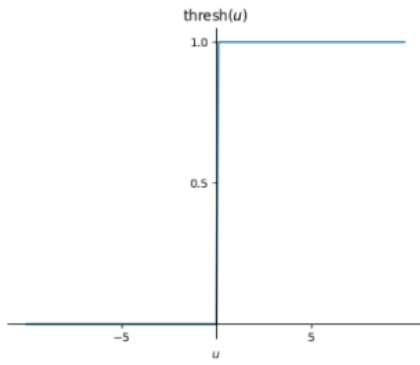
$$\nabla_u J = \frac{\partial J}{\partial u} = \frac{\partial J}{\partial y} = -(d - y) \cdot f'(u)$$

$$\nabla_{\mathbf{w}} J = \frac{\partial J}{\partial \mathbf{w}} = \frac{\partial J}{\partial u} \frac{\partial u}{\partial \mathbf{w}} = -(d - y) \cdot f'(u) \cdot \mathbf{x}$$

$$\nabla_b J = \frac{\partial J}{\partial b} = \frac{\partial J}{\partial u} \frac{\partial u}{\partial b} = -(d - y) \cdot f'(u) \cdot 1$$

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

## Discrete Perceptron



A Discrete Perceptron is given by

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

The activation function is

$$y = f(u) = \begin{cases} 1, & u > 0 \\ 0, & u \leq 0 \end{cases} = \text{thresh}(u)$$

Given a training dataset  $\{(x_p, d_p)\}_{p=1}^P$  where input  $\mathbf{x} \in \mathbf{R}_n$  and target  $d_p \in \{1, 0\}$  training a Discrete Perceptron neuron finds a linear decision boundary in the feature space

Cost function:  $\delta = (d - y) \in \{-1, 0, 1\}$

## Logistic Regression

A logistic regression neuron performs a binary classification of inputs. That is, it classifies inputs into two classes with labels 0 and 1'. The activation of a logistic regression neuron gives the probability of the neuron output belonging to class

Given an input  $\mathbf{x}$ , the activation of the neuron gives  $P(y=1|\mathbf{x})$

$$f(u) = P(y = 1 | \mathbf{x}) = \frac{1}{1 + e^{-u}}$$

Nb: the input gives the probability that the input belong to class

The output  $y$  of the neuron is given by:

$$y = 1(f(u) > 0.5)$$

Decision boundary is given by  $u=0$ .

Cost function:

$$J = -d \cdot \log(f(u)) - (1 - d) \cdot \log(1 - f(u))$$

$$\nabla_u J = \frac{\partial J}{\partial u} = \frac{\partial J}{\partial y} = -(d - f(u)) \cdot f'(u)$$

$$\nabla_{\mathbf{w}} J = \frac{\partial J}{\partial \mathbf{w}} = \frac{\partial J}{\partial u} \frac{\partial u}{\partial \mathbf{w}} = -(d - f(u)) \cdot \mathbf{x}$$

$$\nabla_b J = \frac{\partial J}{\partial b} = \frac{\partial J}{\partial u} \frac{\partial u}{\partial b} = -(d - f(u)) \cdot 1$$

Batch learning	Stochastic learning
$(X, d)$	$(x_p, d_p)$
$\mathbf{u} = \mathbf{X}\mathbf{w} + b\mathbf{1}_p$	$u_p = \mathbf{x}_p^T \mathbf{w} + b$
$\mathbf{y} = 1(\mathbf{u} > 0)$	$y_p = 1(u_p > 0)$
$\mathbf{w} \leftarrow \mathbf{w} + \alpha \mathbf{X}^T (d - \mathbf{y})$	$\mathbf{w} \leftarrow \mathbf{w} + \alpha (d_p - y_p) \mathbf{x}_p$
$b \leftarrow b + \alpha \mathbf{1}_p^T (d - \mathbf{y})$	$b \leftarrow b + \alpha (d_p - y_p)$

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

## Linear separator

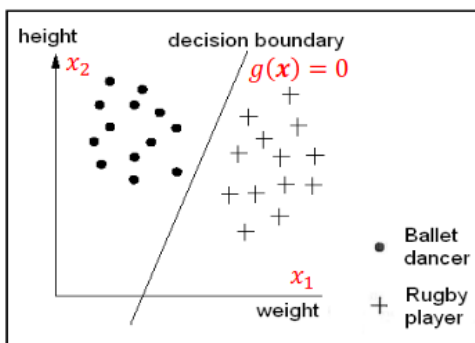


Figure: A linear classification decision boundary

$$G(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

$$\mathbf{w} = (w_1 \ w_2 \ \dots \ w_n), \quad b = w_0$$

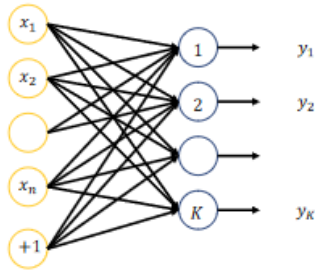
The weight vector is orthogonal to the line.

$$\begin{pmatrix} 1 \\ 3 \end{pmatrix} \perp \begin{pmatrix} 3 \\ 1 \end{pmatrix}$$

## Derivatives

	$y$	$y'$
<b>Linear</b>	$u$	1
<b>Perceptron "Sigmoid"</b>	$\frac{1}{1 + e^u}$	$y(1 - y)$
<b>Tanh "Bipolar Sigmoid"</b>	$\frac{e^u - e^{-u}}{e^u + e^{-u}}$	$1 - y^2$
<b>ReLU</b>	$\max(0, u)$	NA
<b>Threshold</b>	$1(u > 0)$	NA

## Layers – Single input



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

Weights connected a neuron layers is represented by a weight  $\mathbf{W}$  where columns are given by weight vectors connected to individual neurons. The dimension is (feature X labels)

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

$$= \text{features} \left\{ \begin{array}{cccc} w_{11} & w_{12} & \dots & w_{1K} \\ w_{21} & w_{22} & \dots & w_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n1} & w_{n2} & \dots & w_{nK} \end{array} \right\}$$

*labels*

And

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

Given an input pattern  $\mathbf{x} \in \mathbf{R}^n$  to a layer of  $k$  neurons  
 $\phi: \mathbf{R}^n \rightarrow \mathbf{R}^K$

Synaptic input to  $k$ th neuron  $u_k$ :

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

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

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

## Layers – Batch

Given a training dataset  $\{(\mathbf{x}_p)\}_{p=1}^P$  input patterns to a layer of  $K$  neurons where  $\mathbf{x}_p \in \mathbf{R}^n$

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 of  $P$  pattern

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

With bias matrix

$$\mathbf{B} = \begin{pmatrix} \mathbf{b}^T \\ \mathbf{b}^T \\ \vdots \\ \mathbf{b}^T \end{pmatrix}$$

And input with dimension (pattern X features)

$$\mathbf{X} = \text{patterns} \left\{ \begin{array}{cccc} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{P1} & x_{P2} & \dots & x_{Pn} \end{array} \right\}$$

*Features*

And activation function

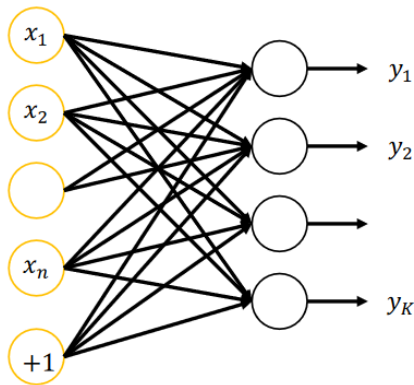
$$\mathbf{Y} = F(\mathbf{U}) = \begin{pmatrix} f(\mathbf{u}_1)^T \\ f(\mathbf{u}_2)^T \\ \vdots \\ f(\mathbf{u}_P)^T \end{pmatrix} = \begin{pmatrix} y_1^T \\ y_2^T \\ \vdots \\ y_P^T \end{pmatrix}$$

$$= P_{\text{pattern}} \left\{ \begin{array}{cccc} y_{11} & y_{12} & \dots & y_{1K} \\ y_{21} & y_{22} & \dots & y_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ y_{P1} & y_{P2} & \dots & y_{PK} \end{array} \right\}$$

*Kneurons*

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

## Perceptron layer



A layer that perform multidimensional regression. A Layer of K perceptron learns a multidimensional non-linear mapping:

$$\phi: R^n \rightarrow R^K$$

The learning is like single neuron learning.

Cost function:  $J = \frac{1}{2} (d - y)^2$

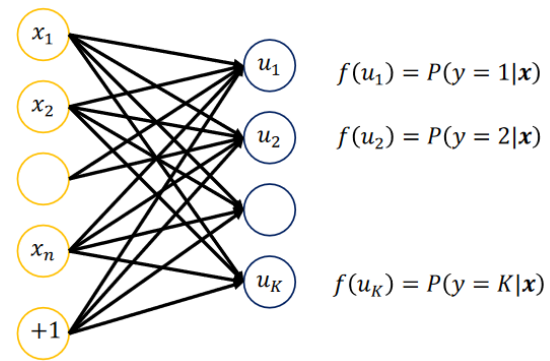
$$\nabla_{\mathbf{u}} J = \begin{pmatrix} \nabla_{u_1} J \\ \nabla_{u_2} J \\ \vdots \\ \nabla_{u_K} J \end{pmatrix} = -(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})$$

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

## Summary for layers

Layers	$f(\mathbf{U}), \mathbf{Y}$	$\nabla_U J$
Linear	$\mathbf{Y} = f(\mathbf{U}) = \mathbf{U}$	$-(\mathbf{D} - \mathbf{Y})$
Perceptron	$Y = f \frac{1}{1 + e^{-U}}$	$-(\mathbf{D} - \mathbf{Y}) \cdot f'(\mathbf{U})$
Softmax	$F(\mathbf{U}) = \frac{e^U}{\sum_{k=1}^K e^{U_k}}$ $\mathbf{y} = \underset{k}{\operatorname{argmax}} f(\mathbf{U})$	$-(K - f(\mathbf{U}))$ K= Onehot

## Softmax Layer



Softmax layer is the extension of a logistic regression to multiclass classification problem, which is also know as multinomial logistic regression

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 x belongs to 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 the weight vector and  $b_k$  is bias of neuron k

The above activation function f is known as softmax activation function

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

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

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

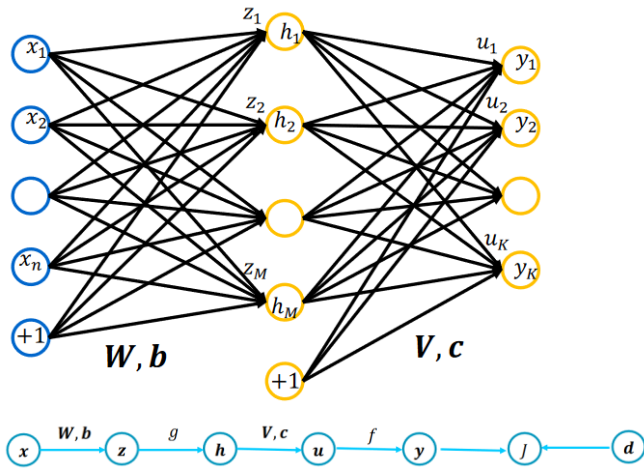
$$\mathbf{u} = \begin{pmatrix} 0.86 & 0.11 & 0.64 \\ -0.84 & -0.28 & -0.49 \\ -0.41 & -0.16 & -0.22 \\ -0.01 & -0.21 & 0.17 \\ -0.86 & -0.82 & -0.06 \\ -1.15 & -0.27 & -0.75 \\ 0.11 & -0.31 & 0.35 \\ -0.12 & -0.10 & -0.02 \end{pmatrix}$$

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

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

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

## Deep FFN - SGD



Synaptic input  $\mathbf{z}$  to hidden layer:

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

Output  $\mathbf{h}$  of the hidden layers

$$\mathbf{h} = g(\mathbf{z})$$

Synaptic input  $\mathbf{u}$  to output layer:

$$\mathbf{u} = \mathbf{V}^T \mathbf{h} + \mathbf{c}$$

Output  $\mathbf{y}$  of the output layer

$$\mathbf{y} = f(\mathbf{u})$$

Output layer:

$$\nabla_{\mathbf{u}} J = \begin{cases} -(\mathbf{d} - \mathbf{y}), & \text{for linear layers} \\ -(1(\mathbf{k} = \mathbf{d}) - f(\mathbf{u})), & \text{for softmax layer} \end{cases}$$

Update of weight:

$$\begin{aligned} \mathbf{V} &\leftarrow \mathbf{V} - \alpha \cdot \mathbf{h} \cdot (\nabla_{\mathbf{u}} J)^T \\ \mathbf{c} &\leftarrow \mathbf{c} - \alpha \cdot \nabla_{\mathbf{u}} J \end{aligned}$$

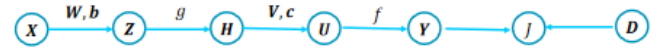
Hidden layer:

$$\nabla_{\mathbf{z}} J = \mathbf{V} \cdot \nabla_{\mathbf{u}} J \cdot g'(\mathbf{z})$$

Update of weight:

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

## Deep FFN – GD



Synaptic input  $\mathbf{Z}$  to hidden layer:

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

Output  $\mathbf{h}$  of the hidden layers

$$\mathbf{H} = g(\mathbf{Z})$$

Synaptic input  $\mathbf{u}$  to output layer:

$$\mathbf{U} = \mathbf{HV} + \mathbf{C}$$

Output  $\mathbf{y}$  of the output layer

$$\mathbf{Y} = f(\mathbf{U})$$

Output layer:

$$\nabla_{\mathbf{u}} J = \begin{cases} -(\mathbf{D} - \mathbf{Y}), & \text{for linear layers} \\ -(K - f(\mathbf{U})), & \text{for softmax layer} \end{cases}$$

Update of weight:

$$\begin{aligned} \mathbf{V} &\leftarrow \mathbf{V} - \alpha \cdot \mathbf{H}^T \cdot (\nabla_{\mathbf{u}} J)^T \\ \mathbf{c} &\leftarrow \mathbf{c} - \alpha \cdot (\nabla_{\mathbf{u}} J)^T \cdot \mathbf{1}_p \end{aligned}$$

Hidden layer:

$$\nabla_{\mathbf{z}} J = \mathbf{V} \cdot \nabla_{\mathbf{u}} J \cdot g'(\mathbf{Z})$$

Update of weight:

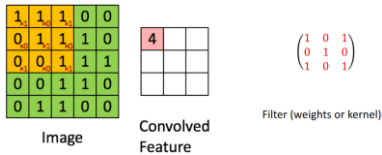
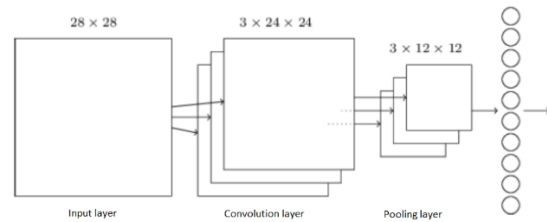
$$\begin{aligned} \mathbf{W} &\leftarrow \mathbf{W} - \alpha \cdot \mathbf{X}^T \cdot \nabla_{\mathbf{z}} J \\ \mathbf{b} &\leftarrow \mathbf{b} - \alpha \cdot (\nabla_{\mathbf{z}} J)^T \cdot \mathbf{1}_p \end{aligned}$$

Definition

	SGD	GD
Single Neuron	$u_p = \mathbf{w}^T \mathbf{x}_p + b$ $y_p = f(u_p)$	$\mathbf{u} = \mathbf{X}\mathbf{w} + b$ $\mathbf{y} = f(\mathbf{u})$
Layered	$\mathbf{u} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$ $\mathbf{y} = f(\mathbf{u})$	$\mathbf{U} = \mathbf{X}\mathbf{W} + \mathbf{B}$ $\mathbf{Y} = f(\mathbf{U})$

$\mathbf{x} = \text{features} \left\{ \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \right.$	$\mathbf{X} = \text{patterns} \left\{ \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_p^T \end{pmatrix} \right. = \text{patterns} \left\{ \underbrace{\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{p1} & x_{p2} & \dots & x_{pn} \end{bmatrix}}_{\text{Features}} \right.$
$\mathbf{w} = \text{features} \left\{ \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \right.$	$\mathbf{W} = \underbrace{(\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_K)}_{\text{labels}} = \text{features} \left\{ \underbrace{\begin{bmatrix} w_{11} & w_{12} & \dots & w_{1K} \\ w_{21} & w_{22} & \dots & w_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n1} & w_{n2} & \dots & w_{nK} \end{bmatrix}}_{\text{labels}} \right.$
$\mathbf{b} = \text{labels} \left\{ \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_K \end{bmatrix} \right.$	$\mathbf{B} = \text{patterns} \left\{ \begin{bmatrix} \mathbf{b}^T \\ \mathbf{b}^T \\ \vdots \\ \mathbf{b}^T \end{bmatrix} \right. = \text{patterns} \left\{ \underbrace{\begin{bmatrix} b_1 & b_2 & \dots & b_K \\ b_1 & b_2 & \dots & b_K \\ \vdots & \vdots & \ddots & \vdots \\ b_1 & b_2 & \dots & b_K \end{bmatrix}}_{\text{labels}} \right.$
$\mathbf{y} = \text{labels} \left\{ \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_K \end{bmatrix} \right.$	$\mathbf{Y} = \begin{pmatrix} y_1^T \\ y_2^T \\ \vdots \\ y_P^T \end{pmatrix} = \text{patterns} \left\{ \underbrace{\begin{bmatrix} y_{11} & y_{12} & \dots & y_{1K} \\ y_{21} & y_{22} & \dots & y_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ y_{P1} & y_{P2} & \dots & y_{PK} \end{bmatrix}}_{K_{\text{neurons}}} \right.$

## CNN



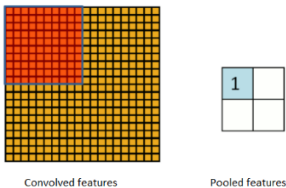
Let  $x_p$  denote the input window of the pixel, and  $\mathbf{w}$  be the filter of the feature map, then the output of the neuron will be

$$u_p = \mathbf{x}_p^T \mathbf{w} + b$$

$$y_p = \text{sigmoid}(u_p)$$

3 filters give 3 feature maps.

When the convolution features are found the pooling is done by taking "max" or "avg" pooling.



### Size of feature maps

$I \times J$  = image\_input size  $L \times M$  = window size

1. VALID: apply the filter wherever it completely overlaps with the input.
  - $(I - L + 1) \times (J - M + 1)$
2. SAME: apply the filter to make the output same size as input. The input is padded with zeros in obtaining the output.
  - $I \times J$



### Stride:

The distance between adjacent centers of the kernel:

The convolution operation is intended to extract features of the input space as finely as possible. The **default** strides for convolution is 1 unit.

Pooling is intended to subsample the convolution layer. The **default** stride for pooling is equal the filter width

## CNN – Example

Given an input pattern  $\mathbf{X}$ :

$$\mathbf{X} = \begin{pmatrix} 0.5 & -0.1 & 0.2 & 0.3 & 0.5 \\ 0.8 & 0.1 & -0.5 & 0.5 & 0.1 \\ -1.0 & 0.2 & 0.0 & 0.3 & -0.2 \\ 0.7 & 0.1 & 0.2 & -0.6 & 0.3 \\ -0.4 & 0.0 & 0.2 & 0.3 & -0.3 \end{pmatrix}$$

The input pattern is received by a convolution layer consisting of one kernel

$$\mathbf{w} = \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix}, \quad b = 0.05$$

If convolution layer has a sigmoid activation function, find the outputs of the convolution layer if the padding is VALID at strides = 1.

If the pooling layer uses max pooling, has a pooling window size of 2x2, and strides = 2, find the activations at the pooling layer for VALID and SAME padding.

Synaptic input to the convolution layer

$$f(u(i, j)) = \sum_l \sum_m x(i + l, j + m) \cdot w(l, m) + b$$

$$\mathbf{U} = \begin{pmatrix} -0.35 & 1.35 & 0.75 \\ -0.55 & 0.85 & 0.05 \\ 0.75 & 0.05 & 1.15 \end{pmatrix}$$

Output of the convolution layer

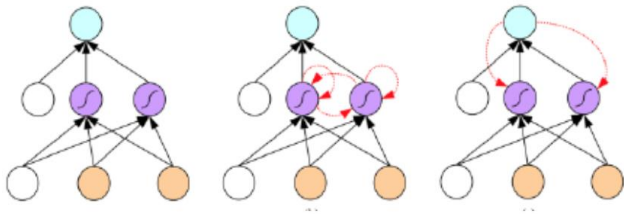
$$f(\mathbf{U}) = \frac{1}{1 + e^{-\mathbf{U}}} = \begin{pmatrix} 0.413 & 0.794 & 0.679 \\ 0.366 & 0.701 & 0.512 \\ 0.679 & 0.512 & 0.76 \end{pmatrix}$$

Output of the max pooling layer:

$$(0.794) : \text{VALID pooling}$$

$$\begin{pmatrix} 0.794 & 0.670 \\ 0.679 & 0.76 \end{pmatrix} : \text{SAME pooling}$$

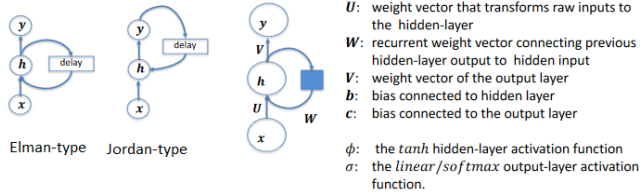
## RNN



Feedforward NN

RNN with hidden recurrence (Elman-type)

RNN with top-down recurrence (Jordan-type)



## SGD

Let  $\mathbf{x}(t)$ ,  $\mathbf{y}(t)$  and  $\mathbf{h}(t)$  be the input, output and hidden output of the network at time  $t$ .

Activation of the Elman-type RNN with one hidden-layer is given by:

$$\mathbf{h}(t) = \Phi(\mathbf{U}^T \mathbf{x}(t) + \mathbf{W}^T \mathbf{h}(t-1) + \mathbf{b})$$

$$\mathbf{y}(t) = \sigma(\mathbf{V}^T \mathbf{h}(t) + \mathbf{c})$$

## GD

Given  $P$  patterns  $\{\mathbf{x}_p\}_{p=1}^P$  where

$$\mathbf{x}_p = (\mathbf{x}_p(t))_{t=1}^T$$

denoting time from 1 to  $T$

$$\mathbf{X}(t) = \begin{pmatrix} \mathbf{x}_1(t)^T \\ \mathbf{x}_2(t)^T \\ \vdots \\ \mathbf{x}_p(t)^T \end{pmatrix}$$

Let  $\mathbf{X}(t)$ ,  $\mathbf{Y}(t)$  and  $\mathbf{H}(t)$  be batch input, output and hidden output of the network at time  $t$ .

Activation of the three-layer Elman-type RNN is given by:

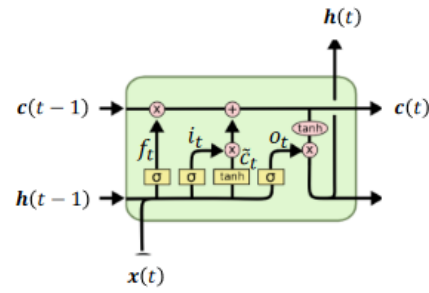
$$\mathbf{H}(t) = \phi \left( \underbrace{\mathbf{X}(t)\mathbf{U}}_{\text{input} \cdot \mathbf{W}_1} + \underbrace{\mathbf{H}(t-1)\mathbf{W}}_{\text{Hidden out} \cdot \mathbf{W}_2} + \mathbf{B} \right)$$

$$\mathbf{Y}(t) = \sigma(\mathbf{H}(t)\mathbf{V} + \mathbf{C})$$

## Matrix sizes

RNN	LSTM	GRU
$\mathbf{w}$	$\{\mathbf{w}_i \ \mathbf{w}_f \ \mathbf{w}_o \ \mathbf{w}_c\}$	$\{\mathbf{w}_r \ \mathbf{w}_z \ \mathbf{w}_h\}$
$\mathbf{b}$	$\{\mathbf{b}_i \ \mathbf{b}_f \ \mathbf{b}_o \ \mathbf{b}_c\}$	$\{\mathbf{b}_r \ \mathbf{b}_z \ \mathbf{b}_h\}$
$\mathbf{U}$	$\{\mathbf{U}_i \ \mathbf{U}_f \ \mathbf{U}_o \ \mathbf{U}_c\}$	$\{\mathbf{U}_r \ \mathbf{U}_z \ \mathbf{U}_h\}$

## Long Short Term Memory



**F is the forget gate:**

$$\mathbf{f}(t) = \sigma(\mathbf{U}_f^T \mathbf{x}(t) + \mathbf{W}_f^T \mathbf{h}(t-1) + \mathbf{b}_f)$$

The forget gate can modulate the memory cell's self-recurrent connection, allowing the cell to remember or forget its previous state as needed. The value of  $\mathbf{f}(t)$  determines if  $\mathbf{c}(t-1)$  is to be remembered or not.

**I is the input gate,  $\tilde{c}$  is the internal cell state**

$$\mathbf{i}(t) = \sigma(\mathbf{U}_i^T \mathbf{x}(t) + \mathbf{W}_i^T \mathbf{h}(t-1) + \mathbf{b}_i)$$

$$\tilde{\mathbf{c}}(t) = \phi(\mathbf{U}_c^T \mathbf{x}(t) + \mathbf{W}_c^T \mathbf{h}(t-1) + \mathbf{b}_c)$$

The input gate can allow incoming signal to alter the state of the memory cell or block it. It decides what new information to store in the cell stage

This has two parts: a sigmoid input gate layer decides which values to update; a  $\tanh$  layer creates a vector of new candidate values  $\tilde{\mathbf{c}}$  that could be added to the state.

**C is the cell state**

$$\mathbf{c}(t) = \tilde{\mathbf{c}}(t) \odot \mathbf{i}(t) + \mathbf{c}(t-1) \odot \mathbf{f}(t)$$

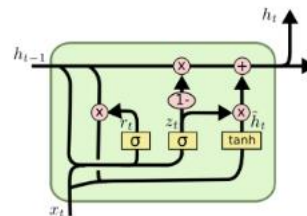
**O is the output gate**

$$\mathbf{o}(t) = \sigma(\mathbf{U}_o^T \mathbf{x}(t) + \mathbf{W}_o^T \mathbf{h}(t-1) + \mathbf{b}_o)$$

$$\mathbf{h}(t) = \phi(\mathbf{c}(t)) \odot \mathbf{i}(t)$$

The output gate can allow the state of the memory cell to have an effect on other neurons or prevent it

## Gated Recurrent Unit



$$\mathbf{r}(t) = \sigma(\mathbf{U}_r^T \cdot \mathbf{x}(t) + \mathbf{W}_r^T \cdot \mathbf{h}(t-1) + \mathbf{b}_r)$$

$$\mathbf{z}(t) = \sigma(\mathbf{U}_z^T \cdot \mathbf{x}(t) + \mathbf{W}_z^T \cdot \mathbf{h}(t-1) + \mathbf{b}_z)$$

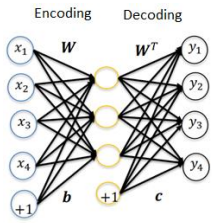
$$\tilde{\mathbf{h}}(t) = \phi(\mathbf{U}_h^T \cdot \mathbf{x}(t) + \mathbf{W}_h^T (r(t) \odot \mathbf{h}(t-1)) + \mathbf{b}_h)$$

$$\mathbf{h}(t) = (1 - \mathbf{z}(t)) \odot \mathbf{h}(t-1) + \mathbf{z}(t) \odot \tilde{\mathbf{h}}(t)$$

$\mathbf{r}$  is the reset gate and  $\mathbf{z}$  update gate.



## Autoencoders



The input signal is encoded into a code that represent the input, the output decodes the signal again. The autoencoder goal is to reproduce the input data. The network therefore learn the features of the input.

The network only performs optimal on trained data variation of trained inputs.

Given an input  $\mathbf{x}$ , the hidden-layer performs the encoding function  $\mathbf{h} = \phi(\mathbf{x})$  and the decoder  $\varphi$  produces the reconstruction  $\mathbf{y} = \varphi(\mathbf{h})$ .

If the autoencoder succeeds:

$$\mathbf{y} = \varphi(\phi(\mathbf{x})) = \mathbf{x}$$

Reverse mapping form the hidden layer to the output can be optionally constrained to be the same as the input to hidden layer. That is  $\mathbf{W}' = \mathbf{W}^T$

Hidden layer and output layer activation can be written as

$$\begin{aligned}\mathbf{h} &= g(\mathbf{W}^T \mathbf{x} + \mathbf{b}) : g(\mathbf{XW} + \mathbf{B}) \\ \mathbf{y} &= f(\mathbf{W}\mathbf{h} + \mathbf{c}) : f(\mathbf{HW}^T + \mathbf{B}')\end{aligned}$$

$f$  is usually sigmoid

The cost function of reconstruction can be measured by many ways, depending on the appropriate distributional assumptions on the inputs.

Learning of autoencoders is unsupervised as no specific targets are given.

Given a training set  $\{\mathbf{x}_p\}_{p=1}^P$

The mean-square-error cost is usually used if the data is assumed to be continuous and Gaussian distributed:

$$J_{mse} = \frac{1}{P} \sum_{p=1}^P |\mathbf{y}_p - \mathbf{x}_p|^2$$

Where  $\mathbf{y}_p$  is the output for input  $\mathbf{x}_p$  and  $|\cdot|$  denotes the magnitude of the vector.

If the inputs are interpreted as bit vectors or vectors of bit probabilities, cross-entropy of the reconstruction can be used

$$J_{cross} = - \sum_{p=1}^P (\mathbf{x}_p \cdot \log \mathbf{y}_p + (1 - \mathbf{x}_p) \cdot \log(1 - \mathbf{y}_p))$$

## Denoising Autoencoders (DAE)

This encoder is used to removed noise from the data

The denoising autoencoders (DAE) receives corrupted data points as inputs and is trained to predict the original uncorrupted data points as its output.

The idea of DAE is that in order to force the hidden layer to discover more robust features and prevent it from simply learning the identity. We train the autoencoder to reconstruct the input from a corrupted version of it.

In other words, DAE attempts to encode the input (preserve the information about input) and attempts to undo the effect of corruption process applied to the input of the autoencoder.

To obtain corrupted version of input data, each input  $x_i$  of input data is added with additive or multiplicative noise.

Additive noise:

$$\tilde{x}_i = x_i + \epsilon$$

where noise  $\epsilon$  is Gaussian distributed:

$$\epsilon \sim N(0, \sigma^2)$$

And  $\sigma$  is the standard deviation that determines the ration S/N ratio. Usually used for continuous data.

Multiplication noise:

$$\tilde{x}_i = \epsilon \cdot x_i$$

where noise  $\epsilon$  could be Binomially distributed:

$$\epsilon \sim \text{Binomial}(p)$$

And  $p$  is the probability of ones and  $1 - p$  is the probability of zeros (noise). Usually, used for binary data

Reconstruction error

$$\text{Error} = \mathbf{X} - \hat{\mathbf{X}} = \begin{pmatrix} e_{11}^T \\ e_{21}^T \end{pmatrix} = n \begin{pmatrix} e_{11} & e_{12} & \dots & e_{1K} \\ e_{21} & e_{22} & \dots & e_{2K} \end{pmatrix}$$

Under and overcomplete

There are two types of autoencoders, undercomplete and overcomplete

Input dimension  $n$  and hidden dimension  $M$ :

if  $M < n$ , undercomplete

if  $M > n$ , overcomplete

Undercomplete reduce the information given, and it therefore need to learn the most important features of the data.

The overcomplete need to be forced to learn the important features by a regulator

## Sparse Autoencoders (SAE)

A SAE is simply an autoencoder whose training criterion involves the sparsity penalty  $\Omega_{sparsity}(\mathbf{h})$  at the hidden layer

$$J_1 = J + \beta_{sparsity}(\mathbf{h})$$

The sparsity penalty term makes the features (weights) learnt by the hidden-layer to be sparse

With sparsity constraint, one would constraint the neurons at the hidden layers to be inactive for most of the time. We say that the neuron is "active" when its output is close to 1 and the neuron is "inactive" when its output is close to 0

For a set  $\{\mathbf{x}_p\}_{p=1}^P$  of input patterns, the average activation  $\rho_j$  of neuron j at the hidden-layer is given by

$$\rho_j = \frac{1}{P} \sum_{p=1}^P h_{pj} = \frac{1}{P} \sum_{p=1}^P f(\mathbf{x}_p^T \cdot \mathbf{w}_j + b)$$

Where  $h_{pj}$  is the activation of hidden neuron j for p th pattern, and  $\mathbf{w}_j$  and  $b_j$  are the weight and biases of the hidden neuron j.

We would like to enforce the constrain  $\rho_j = \rho$  such that the sparsity parameter  $\rho$  is set to a small value close to zero.

That is, most of the time, hidden neuron activations are maintained at  $\rho$  on average. By choosing a smaller value for  $\rho$ , the neurons are activated selectively to patterns and thereby learn sparse features.

To achieve sparse activations at the hidden-layer, the Kullback-Leibler (KL) divergence is used as the sparsity constraint:

$$D(\mathbf{h}) = \sum_{j=1}^M \rho \cdot \log \frac{\rho}{\rho_j} + (1 - \rho) \cdot \log \frac{1 - \rho}{1 - \rho_j}$$

Where M is the number of hidden neurons and  $\rho$  is the sparsity parameter.

KL divergence measures the deviation of the distribution  $\{\rho_j\}$  of activations at the hidden-layer from the uniform distribution  $\rho$ .

The KL divergence is minimum when  $\rho_j = \rho$  for all j. That is, when the average activations are uniform and equal to very low value  $\rho$ .

The penalty is added to the cost function

$$J_1 = J + \beta \cdot D(\mathbf{h})$$

The gradients for can now be found

$$\nabla_{\mathbf{W}} D(\mathbf{h}) = (\nabla_{w_1} D(\mathbf{h}) \quad \nabla_{w_2} D(\mathbf{h}) \quad \dots \quad \nabla_{w_M} D(\mathbf{h}))$$

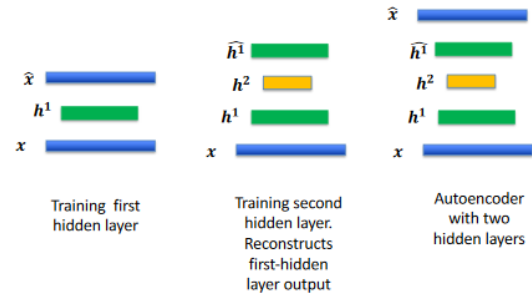
$$\nabla_{\mathbf{b}} D(\mathbf{h}) = (\nabla_{b_1} D(\mathbf{h}) \quad \nabla_{b_2} D(\mathbf{h}) \quad \dots \quad \nabla_{b_M} D(\mathbf{h}))$$

And we update the learning algorithm

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha (\nabla_{\mathbf{W}} J + \beta \cdot \nabla_{\mathbf{W}} D(\mathbf{h}))$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{b}} J + \beta \cdot \nabla_{\mathbf{b}} D(\mathbf{h}))$$

## Deep stacked Autoencoders



Deep autoencoders can be built by stacking autoencoders one after the other. Training of deep autoencoders is done in a step-by-step fashion on layer at time.

After training the first level of denoising autoencoder, the resulting hidden representation is used to train a second level of the denoising encoder. The second level hidden representation can be used to train the third level of the encoders. This process is repeated and deep stacked autoencoder can be realized.

After training a stacked autoencoder, an output layer may be added on the top of the stacked hidden layers for classification/regression. The final layer is a supervised layer such as a Softmax for classification or a linear layer for regression. The parameters of the feedforward network can be fine-tuned to minimize the error in regressing/classifying the targets by supervised gradient descent learning.