

IN THE NAME OF ALLAH

1

Neural Networks

Gradient Descent Learning Rule

Shahrood University of Technology
Hossein Khosravi

Homework #1

2

- Write a program for classification of samples of **two classes** using single layer **perceptron**:
- Samples are in **2D** space
- Sample preparation
 - **Randomly** around a line
 - From 2 **Gaussian** distributions with different μ and σ
- Show the results in each iteration **visually**
- Use **C++**, **MATLAB** or other languages
- **Optional**: User can produce samples interactively
- Repeat the problem with **5 classes**
- Due date: **within 10 days**

Gradient Descent Learning Rule

3

- Gradient Descent: Consider linear unit without threshold and continuous output o (not just $-1,1$)

$$o(X) = w_0 + w_1 x_1 + \dots + w_m x_m = W^T X$$

- The squared error (where D is the set of training examples)
 - ▣ For an example (x,d) the error $e(w)$ of the network is

$$e(W) = d - o(X) = d - \sum_{j=0}^m x_j w_j$$

- ▣ And the squared error of (x,d) is

$$E(W) = \frac{1}{2} e^2$$

Gradient Descent Learning Rule

4

- The total squared error is

$$E(W) = E(w_1, \dots, w_m) = \frac{1}{2} \sum_{(x,d) \in D} (d - o(x))^2$$

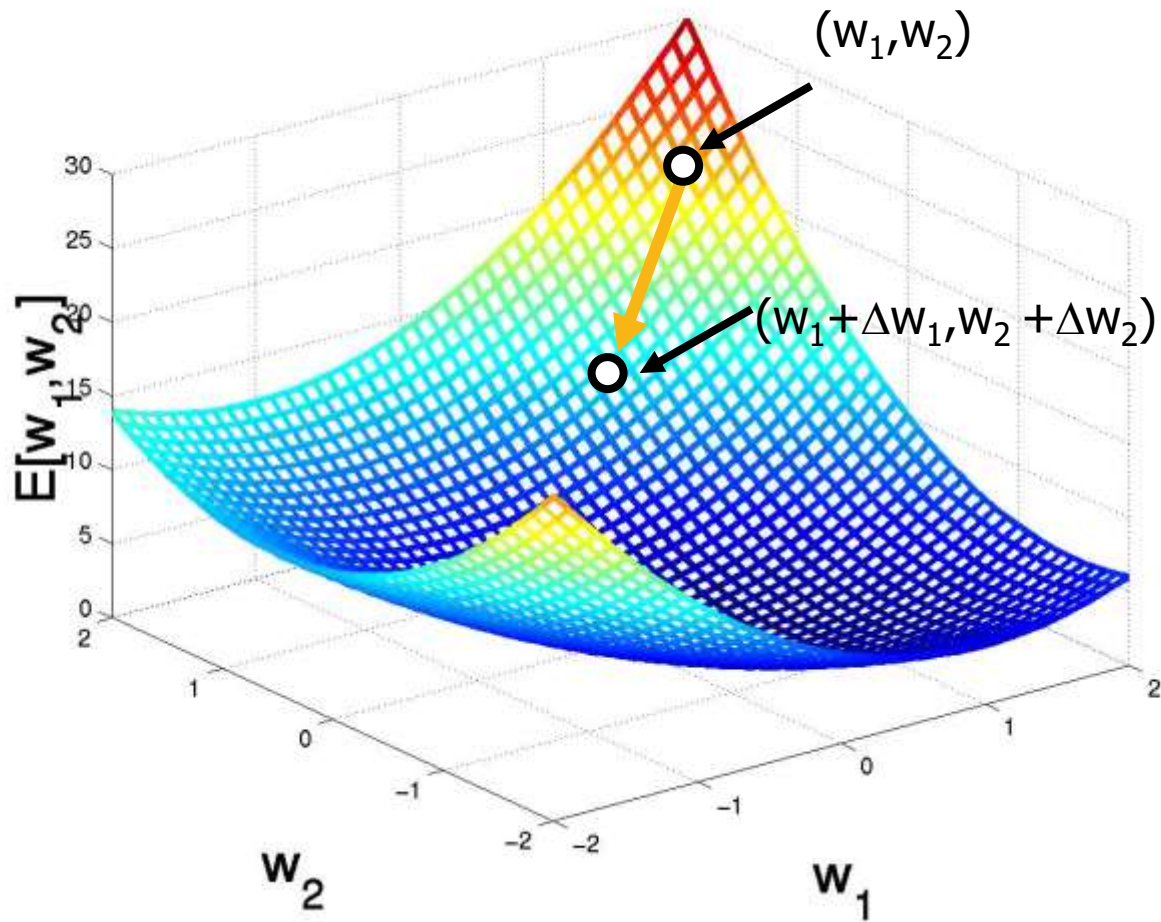
$$E(W) = E(w_1, \dots, w_m) = \frac{1}{2} \sum_{(x,d) \in D} (d - WTX)^2$$

- Update the weight w_i such that $E(W) \rightarrow$ minimum

$$w_i(k + 1) = w_i(k) + \Delta w_i$$

Gradient Descent Learning Rule

5



Gradient Descent Learning Rule

6

- Start from an arbitrary point in the **weight space**
- The direction in which the error E of an example (as a function of the weights) is **decreasing most rapidly** is the opposite of the gradient of E :

$$-\nabla E(\mathbf{W}) = - \left[\frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_m} \right]$$

- Take a small step (of size η) in that direction

$$\mathbf{w}(\mathbf{k} + 1) = \mathbf{w}(\mathbf{k}) - \eta(\nabla E(\mathbf{W}))$$

Gradient Descent Learning Rule

7

- Train the w_i 's such that they minimize the squared error

- $E(W) = E(w_1, \dots, w_m) = \frac{1}{2} \sum_n (d_n - o_n(x))^2$

- Gradient: $\nabla E(w) = \left[\frac{\partial E}{\partial w_0}, \dots, \frac{\partial E}{\partial w_m} \right]$

$$\Delta \mathbf{w} = -\eta \nabla E(\mathbf{w})$$

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} = -\eta \frac{\partial}{\partial w_i} \frac{1}{2} \sum_n (d_n - o_n)^2$$

$$= -\eta \frac{\partial}{\partial w_i} \frac{1}{2} \sum_n (d_n - \sum_i w_i x_{i,n})^2$$

$$= -\eta \sum_n (d_n - o_n) (-x_{i,n})$$

$$= \eta \sum_n (d_n - o_n) x_{i,n}$$

- Gradient descent learning rule

$$w_i(k+1) = w_i(k) + \Delta w_i = w_i(k) + \eta \sum_n (d_n - o_n) x_{i,n}$$

Batch Learning Pseudo Code

Denoting a training example $(x(n), d_n)$ or $((x_{1n}, \dots, x_{mn}), d_n)$ where (x_{1n}, \dots, x_{mn}) are the input values, and d_n is the desired output

- $k=1$, randomly initialize $w_i(k)$, calculate $E(W)$
- **While** ($E(W)$ unsatisfactory && $k < \text{max_iterations}$)
 - Initialize each Δw_i to zero
 - **For** each instance $(x(n), d_n)$ in D do
 - Calculate network output $o_n = \sum_i w_i(k)x_{in}$
 - **For** each weight dimension $w_i(k)$, $i=1, \dots, m$
 - $\Delta w_i = \Delta w_i + \eta(d_n - o_n) x_{in}$
 - **End For**
 - **End For**
 - **For** each weight dimension $w_i(k)$, $i=1, \dots, m$
 - $w_i(k+1) = w_i(k) + \Delta w_i$
 - **EndFor**
 - Calculate $E(W)$ based on updated $w_i(k+1)$
 - $k=k+1$
- **End While**

Example

9

- Consider the 2-dimensional training set $C_1 \cup C_2$,
- $C_1 = \{(1,1), (1, -1), (0, -1)\}$ with class label 1
- $C_2 = \{(-1,-1), (-1,1), (0,1)\}$ with class label 0

- Train a adaline on $C_1 \cup C_2$

Example – small η

- C1: $\{(1, 1, 1), (1, 1, -1), (1, 0, -1)\}$
- C2: $\{(1, -1,-1), (1, -1,1), (1, 0,1)\}$
- $\eta=0.1$
- $w_i(k+1)=w_i(k)+\Delta w_i; \quad \Delta w_i= \eta \sum_n (d_n - o_n) x_{in}$
- Fill out this table sequentially (**First** pass):

Input	$w(k)$	d_n	o_n	$\eta(d_n - o_n)x_{in}$	Δw_i
(1, 1, 1)	(1,0, 0)	1	1	(0, 0, 0)	(0, 0, 0)
(1, 1, -1)	(1, 0, 0)	1	1	(0, 0, 0)	(0, 0, 0)
(1,0, -1)	(1, 0, 0)	1	1	(0, 0, 0)	(0, 0, 0)
(1,-1, -1)	(1, 0, 0)	0	1	(-0.1, 0.1, 0.1)	(-0.1, 0.1, 0.1)
(1,-1, 1)	(1, 0, 0)	0	1	(-0.1, 0.1, -0.1)	(-0.2, 0.2, 0)
(1, 0, 1)	(1, 0, 0)	0	1	(-0.1, 0, -0.1)	(-0.3, 0.2, -0.1)
	E(W)=3/2			$w(k+1)$	(0.7, 0.2, -0.1)

Example – small η

- C1: $\{(1, 1, 1), (1, 1, -1), (1, 0, -1)\}$
- C2: $\{(1, -1,-1), (1, -1,1), (1, 0,1)\}$
- $\eta=0.1$
- $w_i(k+1)=w_i(k)+\Delta w_i;$ $\Delta w_i= \eta \sum_n (d_n - o_n) x_{in}$
- Fill out this table sequentially (**Second** pass):

Input	w(k)	d _n	O _n	$\eta(d_n - o_n)x_{in}$	Δw_i
(1, 1, 1)	(0.7, 0.2, -0.1)	1	0.8	(0.02, 0.02, 0.02)	(0.02,0.02,0.02)
(1, 1, -1)	(0.7, 0.2, -0.1)	1	1	(0,0,0)	(0.02,0.02,0.02)
(1,0, -1)	(0.7, 0.2, -0.1)	1	0.8	(0.02,0,-0.02)	(0.04,0.02,0)
(1,-1, -1)	(0.7, 0.2, -0.1)	0	0.6	(-0.06,0.06,0.06)	(-0.02,0.08,0.06)
(1,-1, 1)	(0.7, 0.2, -0.1)	0	0.4	(-0.04,0.04,-0.04)	(-0.06,0.12,0.02)
(1, 0, 1)	(0.7, 0.2, -0.1)	0	0.6	(-0.06,0,-0.06)	(-0.12,0.12,-0.04)
	E(W)=0.96/2			w(k+1)	(0.58,0.32,-0.14)

Note that unlike perceptron, output can take values other than 0 and 1

Example – small η

12

- C1: $\{(1, 1, 1), (1, 1, -1), (1, 0, -1)\}$
- C2: $\{(1, -1,-1), (1, -1,1), (1, 0,1)\}$
- $\eta=0.1$
- $w_i(k+1)=w_i(k)+\Delta w_i;$ $\Delta w_i= \eta \sum_n (d_n - o_n) x_{in}$
- Fill out this table sequentially (**Third** pass):

Input	w(k)	d _n	O _n	$\eta(d_n - o_n)x_{in}$	Δw_i
(1, 1, 1)	(0.58,0.32,-0.14)	1	0.76	(0.024,0.024,0.024)	(0.024,0.024,0.024)
(1, 1, -1)	(0.58,0.32,-0.14)	1	1.04	(-0.004,-0.004,0.004)	(0.02,0.02,0.028)
(1,0, -1)	(0.58,0.32,-0.14)	1	0.72	(0.028,0,-0.028)	(0.048,0.04,0)
(1,-1, -1)	(0.58,0.32,-0.14)	0	0.4	(-0.06,0.06,0.06)	(-0.012,0.1,0.06)
(1,-1, 1)	(0.58,0.32,-0.14)	0	0.12	(-0.088,0.088,-0.088)	(-0.1,0.188,-0.028)
(1, 0, 1)	(0.58,0.32,-0.14)	0	0.44	(-0.056,0,-0.056)	(-0.156,0.188,-0.084)
	E(W)=0.5056/2			w(k+1)	(0.424,0.508,-0.224)

Gradient Descent Learning Rule

13

- Because the error surface contains only a **single global minimum**, this algorithm will converge to a weight vector with minimum error, regardless of whether the training examples are linearly separable, given a sufficiently small learning rate η
- If η is too large
 - ▣ Overstepping the minimum in the error surface
- Gradually reduce the value of η as the number of gradient descent steps grows

Example – large η

14

- C1: $\{(1, 1, 1), (1, 1, -1), (1, 0, -1)\}$
- C2: $\{(1, -1,-1), (1, -1,1), (1, 0,1)\}$
- $\eta=1$
- $w_i(k+1)=w_i(k)+\Delta w_i;$ $\Delta w_i= \eta \sum_n (d_n - o_n) x_{in}$
- Fill out this table sequentially (**First** pass):

Input	w(k)	d	O	$\eta(d-o)x_i$	Δw_i
(1, 1, 1)	(1,0, 0)	1	1	(0, 0, 0)	(0, 0, 0)
(1, 1, -1)	(1, 0, 0)	1	1	(0, 0, 0)	(0, 0, 0)
(1,0, -1)	(1, 0, 0)	1	1	(0, 0, 0)	(0, 0, 0)
(1,-1, -1)	(1, 0, 0)	0	1	(-1, 1, 1)	(-1, 1, 1)
(1,-1, 1)	(1, 0, 0)	0	1	(-1, 1, -1)	(-2, 2, 0)
(1, 0, 1)	(1, 0, 0)	0	1	(-1, 0, -1)	(-3, 2, -1)
	$E(W)=3/2$			w(k+1)	(-2, 2, -1)

Example – large η

15

- C1: $\{(1, 1, 1), (1, 1, -1), (1, 0, -1)\}$
- C2: $\{(1, -1,-1), (1, -1,1), (1, 0,1)\}$
- $\eta=1$
- $w_i(k+1)=w_i(k)+\Delta w_i;$ $\Delta w_i= \eta \sum_n (d_n - o_n) x_{in}$
- Fill out this table sequentially (**Second** pass):

Input	$w(k)$	d	o	$\eta(d-o)x_i$	Δw_i
(1, 1, 1)	(-2,2, -1)	1	-1	(2, 2, 2)	(2, 2, 2)
(1, 1, -1)	(-2,2, -1)	1	1	(0, 0, 0)	(2, 2, 2)
(1,0, -1)	(-2,2, -1)	1	-1	(2, 0, -2)	(4, 2, 0)
(1,-1, -1)	(-2,2, -1)	0	-3	(3, -3, -3)	(7, -1, -3)
(1,-1, 1)	(-2,2, -1)	0	-5	(5, -5, 5)	(12, -6, 2)
(1, 0, 1)	(-2,2, -1)	0	-3	(3, 0, 3)	(15, -6, 5)
$E(W)=51/2$				$w(k+1)$	(13, -4, 4)

Example – large η

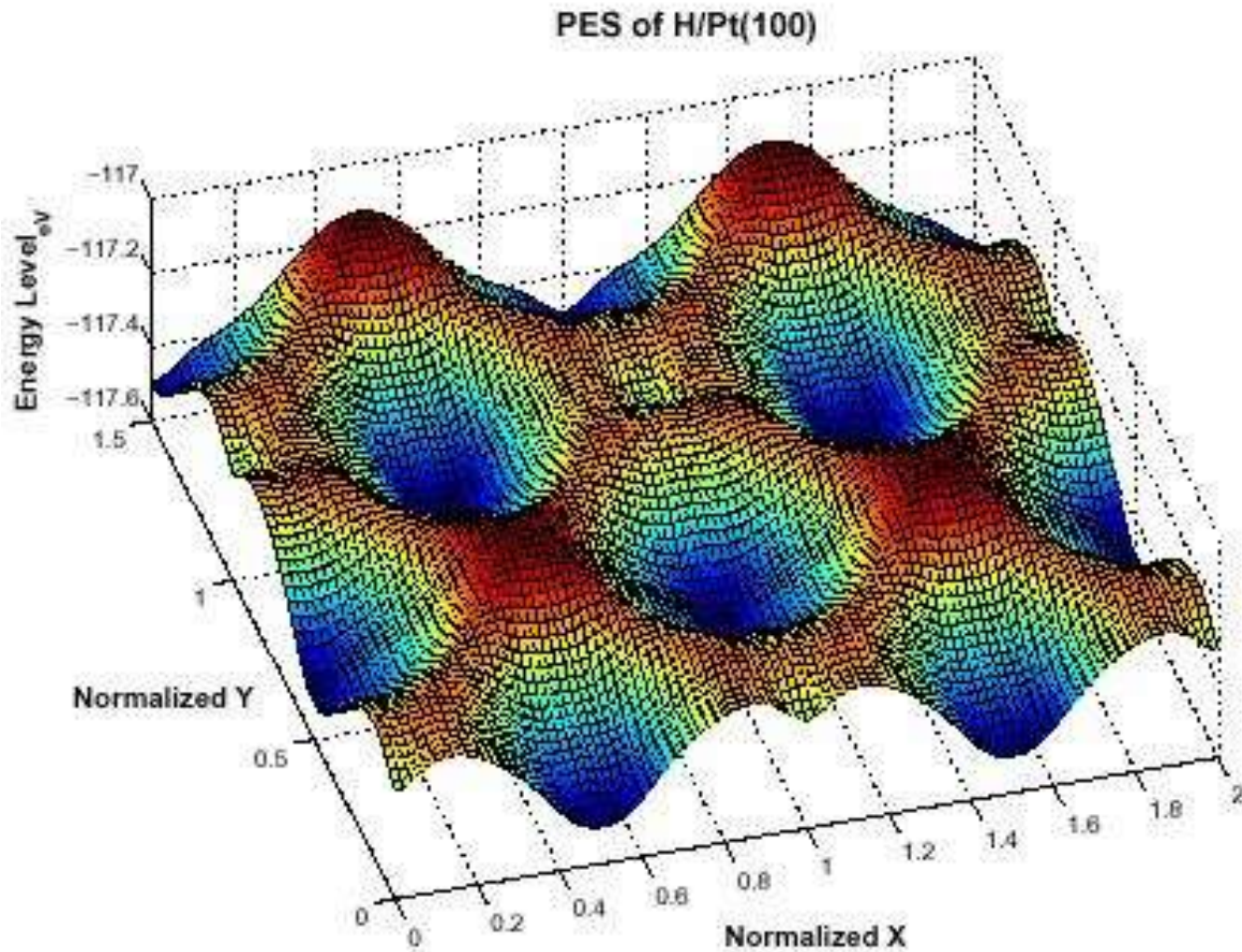
16

- C1: $\{(1, 1, 1), (1, 1, -1), (1, 0, -1)\}$
- C2: $\{(1, -1,-1), (1, -1,1), (1, 0,1)\}$
- $\eta=1$
- $w_i(k+1)=w_i(k)+\Delta w_i;$ $\Delta w_i= \eta \sum_n (d_n - o_n) x_{in}$
- Fill out this table sequentially (**Third** pass):

Input	w(k)	d	O	$\eta(d-o)x_i$	Δw_i
(1, 1, 1)	(13, -4, 4)	1	13	(-12,-12,-12)	
(1, 1, -1)	(13, -4, 4)	1	5	(-4,-4,4)	
(1,0, -1)	(13, -4, 4)	1	9	(-8,0,8)	
(1,-1, -1)	(13, -4, 4)	0	13	(-13,13,13)	
(1,-1, 1)	(13, -4, 4)	0	21	(-21,21,-21)	
(1, 0, 1)	(13, -4, 4)	0	17	(-17,0,-17)	(-75,18,-25)
	$E(W)=1123/2$			$w(k+1)$	(-62,14,-21)

Local minimum

17



Incremental Stochastic Gradient Descent

18

- Batch mode : Gradient descent
 - $w(k+1)=w(k) - \eta \nabla E_D(W)$ **over the entire data D**
 - $E_D(W)=1/2 \sum_n (d_n - o_n)^2$

- Incremental mode: Gradient descent
 - $w(k+1)=w(k) - \eta \nabla E_n(W)$ **over individual training examples**
 - $E_n(W)=1/2 (d_n - o_n)^2$

- Incremental Gradient Descent can approximate Batch Gradient Descent arbitrarily close if η is small enough

Weights Update Rule: incremental mode

19

- Computation of Gradient(E):

$$\begin{aligned}\frac{\partial E(W)}{\partial w} &= \frac{\partial}{\partial w} \left(\frac{1}{2} (d_n - o_n)^2 \right) = \frac{\partial}{\partial w} \left(\frac{1}{2} (d_n - w^T x(n))^2 \right) \\ &= -(d_n - o_n) x(n)\end{aligned}$$

- Delta rule (**AdaLine**: Adaptive Linear Elements) for weight update:

$$\begin{aligned}w(k+1) &= w(k) - \eta \frac{\partial E(W)}{\partial w} \\ w(k+1) &= w(k) + \eta (d_n - o_n) x(n)\end{aligned}$$

Delta Rule (**AdaLine**) learning algorithm

20

k=1;

initialize $w_i(k)$ randomly; Calculate $E_D(W)$

while ($E_D(W)$ unsatisfactory **AND** $k < \text{max_iterations}$)

 Select an example $(x(n), d_n)$

$$\Delta w_i = \eta (d_n - o_n) x_{in}$$

$$w_i(k+1) = w_i(k) + \Delta w_i$$

 Calculate $E_D(W)$

 k = k+1;

end-while;

Example – incremental mode

21

C1: $\{(1, 1, 1), (1, 1, -1), (1, 0, -1)\}$

C2: $\{(1, -1, -1), (1, -1, 1), (1, 0, 1)\}$

$\eta=0.1$

Fill out this table sequentially (First pass):

$$w_i(k+1) = w_i(k) + \eta(d-o)x_i$$

Input	$w_i(k)$	d	O	$\eta(d-o)x_i$	$w_i(k+1)$
(1, 1, 1)	(1,0, 0)	1	1	(0, 0, 0)	(1, 0, 0)
(1, 1, -1)	(1, 0, 0)	1	1	(0, 0, 0)	(1, 0, 0)
(1, 0, -1)	(1, 0, 0)	1	1	(0, 0, 0)	(1, 0, 0)
(1, -1, -1)	(1, 0, 0)	0	1	(-0.1, 0.1, 0.1)	(0.9, 0.1, 0.1)
(1, -1, 1)	(0.9, 0.1, 0.1)	0	0.9	(-0.09, 0.09, -0.09)	(0.81, 0.19, 0.01)
(1, 0, 1)	(0.81, 0.19, 0.01)	0	0.82	(-0.082, 0, -0.082)	(0.728, 0.19, -0.072)

Perceptron Learning Rule vs. Gradient Descent Rule

22

Perceptron learning rule guaranteed to succeed if

- ▣ Training examples are linearly separable
- ▣ Sufficiently small learning rate η

Gradient descent learning rules

- ▣ Guaranteed to converge to hypothesis with minimum squared error
- ▣ Given sufficiently small learning rate η
- ▣ Even when training data contains noise
- ▣ Even when training data not separable by H

Comparison of Perceptron and Adaline

23

	Perceptron	Adaline
Architecture	Single-layer	Single-layer
Neuron model	Non-linear	linear
Learning algorithm	Minimize number of misclassified examples	Minimize total squared error
Application	Linear classification	Linear classification, regression