

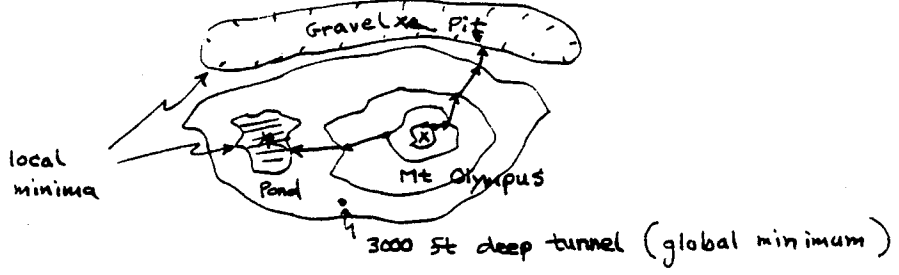
1 May 1988  
Neil E Cottar

# Gradient Descent - Introduction

ex: steepest path = fastest way down

- Water flows in steepest path down a mountain

ex: Given topo map find a local depression (i.e. local minima)



Start at pt on mountain  
Follow gradient = steepest path

- Note:
- 1) Sol'n depends on starting point. • Could end up in pond or in gravel pit
  - 2) Sol'n " " step size in direction of gradient ... • Could step right over the 3000 ft deep mine shaft
  - 3) Sol'n may not be global min

Gradient descent minimizes result (e.g. altitude) with respect to variable parameters (e.g. (attitude, longitude).  
(or  $w_{ij}$ )

## Mathematics:

def: gradient of  $f(x_1, \dots, x_n) \equiv \nabla f \equiv \begin{bmatrix} \partial f / \partial x_1 \\ \vdots \\ \partial f / \partial x_n \end{bmatrix}$

- Multidimensional derivative

ex: multidim Taylor series • See Duda & Hart p.140

~~$f(x_1, \dots, x_n)$~~

$$f(x_1 + \epsilon_1, \dots, x_n + \epsilon_n) = f(x_1, \dots, x_n) + \overset{\text{switch order}}{\nabla f \cdot (\epsilon_1, \dots, \epsilon_n)} + \dots$$

- $\nabla f$  is a function of  $x_1, \dots, x_n$ .
- $\nabla f$  is a vector that has a different value at each point  $(x_1, \dots, x_n)$

Plot  $f(x_1, \dots, x_n)$  as surface in n-dim space.  $\nabla f$  always points in direction of steepest slope. True at every pt  $(x_1, \dots, x_n)$

May 1988  
Neil E Cottler

Gradient Descent - Introduction (cont.)

$$\vec{w}_{new} = \vec{w}_{old} - \rho \nabla f(\vec{w})$$

ex: Perceptron Proportional Increment Training Rule

Let  $f(\vec{v}) = \vec{q} \cdot \vec{v}$   $\vec{v} = (w_0, w_1, \dots, w_n)$   
 $\vec{q} = (1, x_1, \dots, x_n)$

Then  $\nabla f(\vec{v}) = \nabla (-w_0 \cdot 1 - w_1 \cdot x_1 - \dots - w_n \cdot x_n)$

$$= \begin{pmatrix} \frac{\partial}{\partial w_0} ( \quad \quad \quad ) \\ \frac{\partial}{\partial w_1} ( \quad \quad \quad ) \\ \vdots \\ \frac{\partial}{\partial w_n} ( \quad \quad \quad ) \end{pmatrix}$$

$$= \begin{pmatrix} -1 \\ -x_1 \\ \vdots \\ -x_n \end{pmatrix} = -\vec{q}$$

$$\vec{v}_{new} = \vec{v}_{old} - \rho(-\vec{q}) = \vec{v}_{old} + \rho\vec{q}$$

or  $w_{0new} = w_{0old} + \rho$   $\vec{w} = (w_1, \dots, w_n)$   
 $\vec{w}_{new} = \vec{w}_{old} + \rho \vec{x}$   $\vec{x} = (x_1, \dots, x_n)$

In class we used  $\rho=1$

ex: Delta Rules

$$\vec{w}_{new} = \vec{w}_{old} + \rho \underbrace{(V_{out(desired)} - V_{out})}_{\delta} \vec{x}$$

- Is proportional increment training for perceptrons.
- Works on various types of neural networks.