

# Gradient Descent

Sargur Srihari  
[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics

- Simple Gradient Descent/Ascent
- Difficulties with Simple Gradient Descent
- Line Search
- Brent's Method
- Conjugate Gradient Descent
- Weight vectors to minimize error
- Stochastic Gradient Descent

# Gradient-Based Optimization

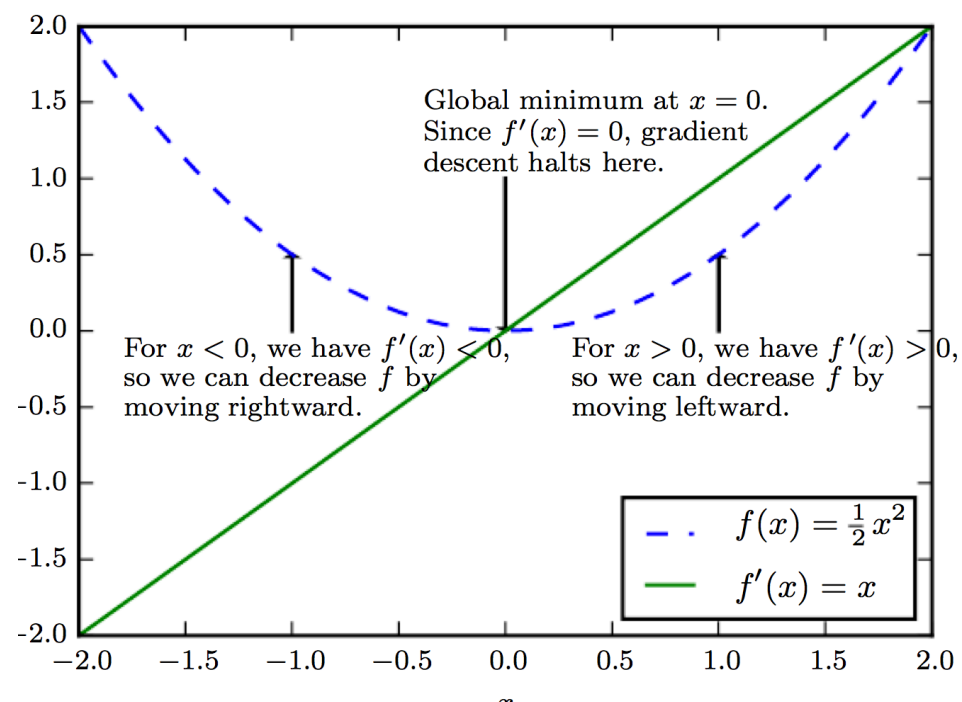
- Most ML algorithms involve optimization
- Minimize/maximize a function  $f(\mathbf{x})$  by altering  $\mathbf{x}$ 
  - Usually stated a minimization
  - Maximization accomplished by minimizing  $-f(\mathbf{x})$
- $f(\mathbf{x})$  referred to as objective function or criterion
  - In minimization also referred to as loss function cost, or error
  - Example is linear least squares  $f(x) = \frac{1}{2} ||Ax - b||^2$
  - Denote optimum value by  $\mathbf{x}^* = \operatorname{argmin} f(\mathbf{x})$

# Calculus in Optimization

- Suppose we have function  $y=f(x)$ ,  $x, y$  real nos.
  - Derivative of function denoted:  $f'(x)$  or as  $dy/dx$ 
    - Derivative  $f'(x)$  gives the slope of  $f(x)$  at point  $x$
    - It specifies how to scale a small change in input to obtain a corresponding change in the output:
$$f(x + \varepsilon) \approx f(x) + \varepsilon f'(x)$$
  - It tells how you make a small change in input to make a small improvement in  $y$
  - We know that  $f(x - \varepsilon \text{ sign}(f'(x)))$  is less than  $f(x)$  for small  $\varepsilon$ . Thus we can reduce  $f(x)$  by moving  $x$  in small steps with opposite sign of derivative
    - This technique is called *gradient descent* (Cauchy 1847)

# Gradient Descent Illustrated

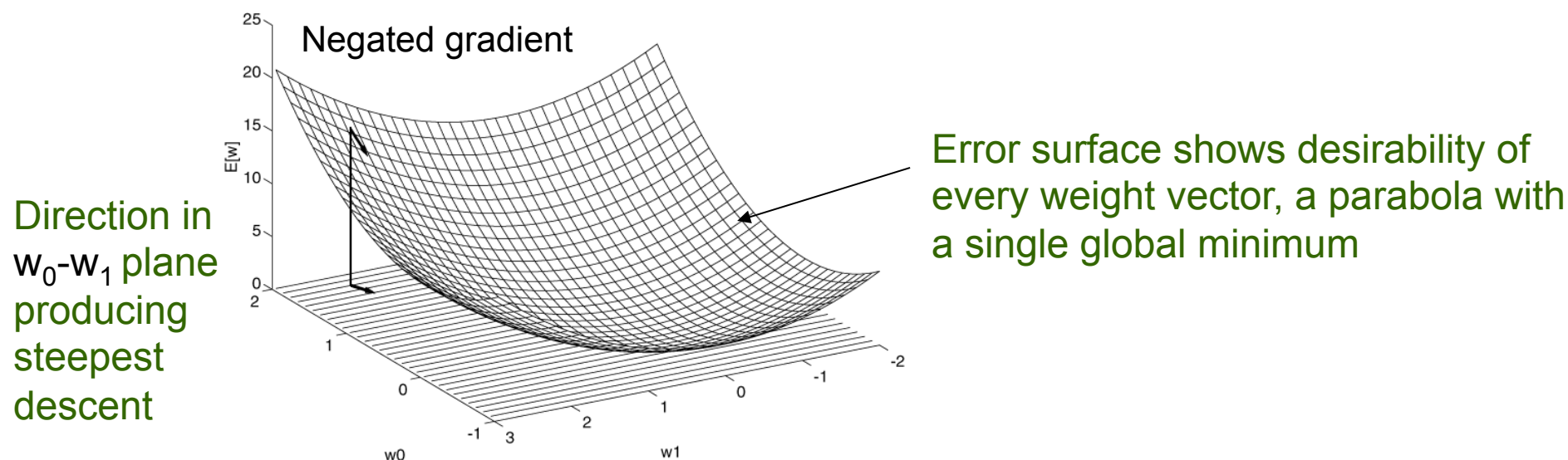
- For  $x > 0$ ,  $f(x)$  increases with  $x$  and  $f'(x) > 0$
- For  $x < 0$ ,  $f(x)$  decreases with  $x$  and  $f'(x) < 0$
- Use  $f'(x)$  to follow function downhill
- Reduce  $f(x)$  by going in direction opposite sign of derivative  $f'(x)$



# Minimizing with multiple inputs

- We often minimize functions with multiple inputs:  $f: R^n \rightarrow R$
- For minimization to make sense there must still be only one (scalar) output

# Application in ML: Minimize Error



- Gradient descent search determines a weight vector  $w$  that minimizes  $E(w)$  by
  - Starting with an arbitrary initial weight vector
  - Repeatedly modifying it in small steps
  - At each step, weight vector is modified in the direction that produces the steepest descent along the error surface

# Definition of Gradient Vector

- The *Gradient (derivative)* of  $E$  with respect to each component of the vector  $w$

$$\nabla E[\vec{w}] \equiv \left[ \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

- Notice  $\nabla E[\vec{w}]$  is a vector of partial derivatives
- Specifies the direction that produces steepest increase in  $E$
- Negative of this vector specifies direction of steepest decrease



# Directional Derivative

- Directional derivative in direction  $\mathbf{u}$  (a unit vector) is the slope of function  $f$  in direction  $\mathbf{u}$

– This evaluates to  $\mathbf{u}^T \nabla_{\mathbf{x}} f(\mathbf{x})$

- To minimize  $f$  find direction in which  $f$  decreases the fastest

– Do this using  $\min_{\mathbf{u}, \mathbf{u}^T \mathbf{u} = 1} \mathbf{u}^T \nabla_{\mathbf{x}} f(\mathbf{x}) = \min_{\mathbf{u}, \mathbf{u}^T \mathbf{u} = 1} \|\mathbf{u}\|_2 \|\nabla_{\mathbf{x}} f(\mathbf{x})\|_2 \cos \theta$

- where  $\theta$  is angle between  $\mathbf{u}$  and the gradient
- Substitute  $\|\mathbf{u}\|_2 = 1$  and ignore factors that not depend on  $\mathbf{u}$  this simplifies to  $\min_{\mathbf{u}} \cos \theta$
- This is minimized when  $\mathbf{u}$  points in direction opposite to gradient
- In other words, the *gradient points directly uphill, and the negative gradient points directly downhill*

# Gradient Descent Rule

$$\vec{w} \leftarrow \vec{w} + \Delta \vec{w}$$

– where

$$\Delta \vec{w} = -\eta \nabla E[\vec{w}]$$

- $\eta$  is a positive constant called the learning rate
  - Determines step size of gradient descent search

- Component Form of Gradient Descent

– Can also be written as

$$w_i \leftarrow w_i + \Delta w_i$$

- where

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

# Method of Gradient Descent

- The gradient points directly uphill, and the negative gradient points directly downhill
- Thus we can decrease  $f$  by moving in the direction of the negative gradient
  - This is known as the method of steepest descent or gradient descent
- Steepest descent proposes a new point

$$\mathbf{x}' = \mathbf{x} - \eta \nabla_{\mathbf{x}} f(\mathbf{x})$$

- where  $\varepsilon$  is the learning rate, a positive scalar. Set to a small constant.

# Simple Gradient Descent

## Procedure Gradient-Descent (

```

 $\theta^1$  //Initial starting point
 $f$  //Function to be minimized
 $\delta$  //Convergence threshold
)
1  $t \leftarrow 1$ 
2 do
3    $\theta^{t+1} \leftarrow \theta^t - \eta \nabla f(\theta^t)$ 
4    $t \leftarrow t + 1$ 
5 while  $\|\theta^t - \theta^{t-1}\| > \delta$ 
6 return  $(\theta^t)$ 

```

## Intuition

Taylor's expansion of function  $f(\theta)$  in the neighborhood of  $\theta^t$  is  $f(\theta) \approx f(\theta^t) + (\theta - \theta^t)^T \nabla f(\theta^t)$

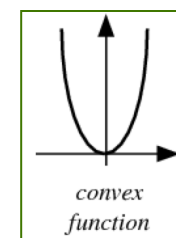
Let  $\theta = \theta^{t+1} = \theta^t + h$ , thus  $f(\theta^{t+1}) \approx f(\theta^t) + h \nabla f(\theta^t)$

Derivative of  $f(\theta^{t+1})$  wrt  $h$  is  $\nabla f(\theta^t)$

At  $h = \nabla f(\theta^t)$  a maximum occurs (since  $h^2$  is positive) and at  $h = -\nabla f(\theta^t)$  a minimum occurs.

Alternatively,

The slope  $\nabla f(\theta^t)$  points to the direction of steepest ascent. If we take a step  $\eta$  in the opposite direction we decrease the value of  $f$



## One-dimensional example

Let  $f(\theta) = \theta^2$

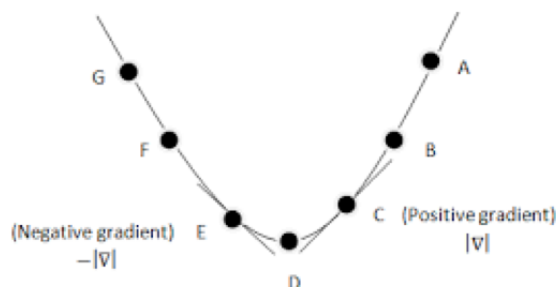
This function has minimum at  $\theta = 0$  which we want to determine using gradient descent

We have  $f'(\theta) = 2\theta$

For gradient descent, we update by  $-f'(\theta)$

If  $\theta^t > 0$  then  $\theta^{t+1} < \theta^t$

If  $\theta^t < 0$  then  $f'(\theta^t) = 2\theta^t$  is negative, thus  $\theta^{t+1} > \theta^t$



# Ex: Gradient Descent on Least Squares

- Criterion to minimize

$$f(\mathbf{x}) = \frac{1}{2} \|A\mathbf{x} - \mathbf{b}\|^2$$

- Least squares regression

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(x_n)\}^2$$

- The gradient is

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = A^T (A\mathbf{x} - \mathbf{b}) = A^T A\mathbf{x} - A^T \mathbf{b}$$

- Gradient Descent algorithm is

1. Set step size  $\varepsilon$ , tolerance  $\delta$  to small, positive nos.

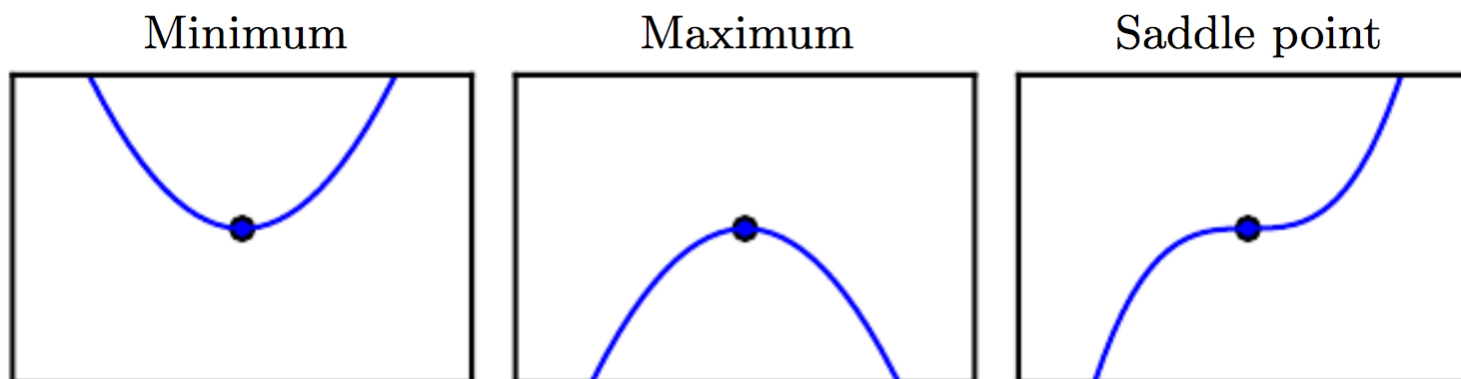
2. *while*  $\|A^T A\mathbf{x} - A^T \mathbf{b}\|_2 > \delta$  *do*

$$\mathbf{x} \leftarrow \mathbf{x} - \eta (A^T A\mathbf{x} - A^T \mathbf{b})$$

3. *end while*

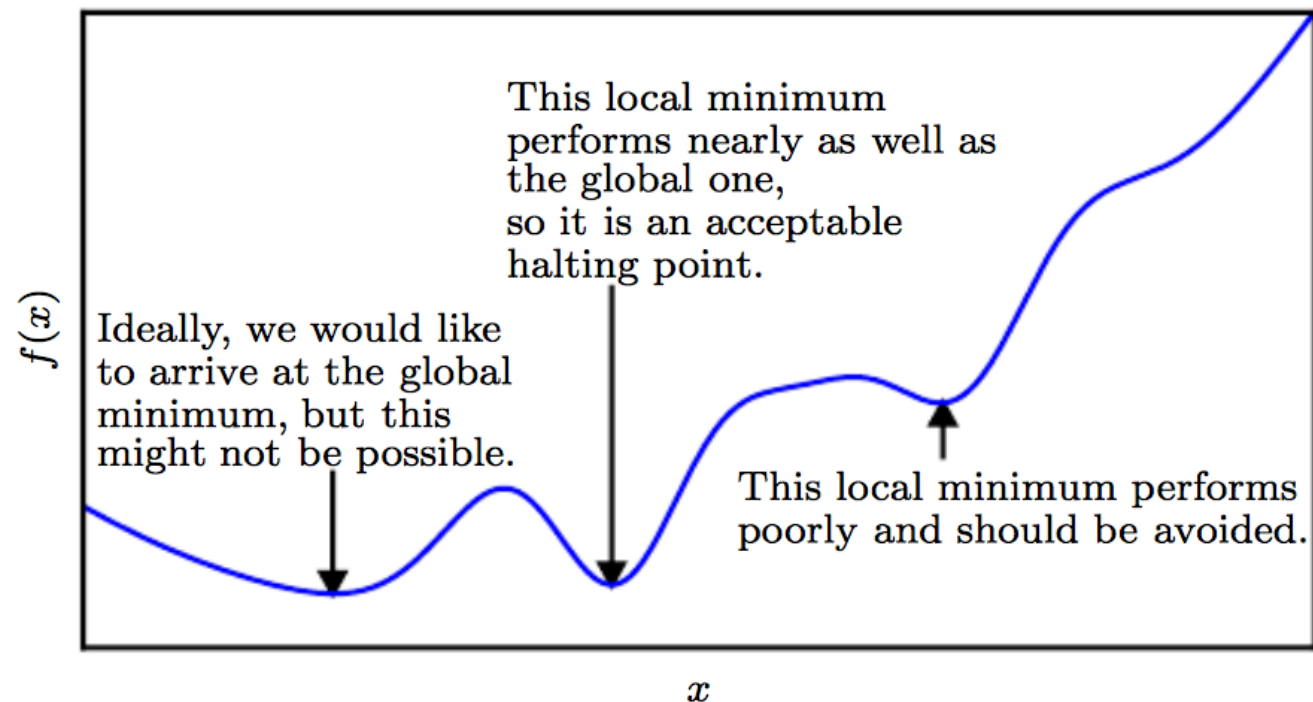
# Stationary points, Local Optima

- When  $f'(x)=0$  derivative provides no information about direction of move
- Points where  $f'(x)=0$  are known as *stationary* or *critical points*
  - Local minimum/maximum: a point where  $f(x)$  lower/higher than all its neighbors
  - Saddle Points: neither maxima nor minima



# Presence of multiple minima

- Optimization algorithms may fail to find global minimum
- Generally accept such solutions



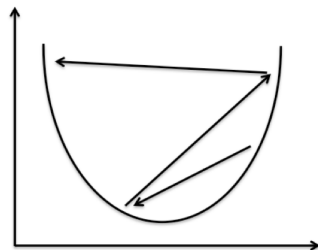
# Difficulties with Simple Gradient Descent

- Performance depends on choice of learning rate  $\eta$

$$\theta^{t+1} \leftarrow \theta^t - \eta \nabla f(\theta^t)$$

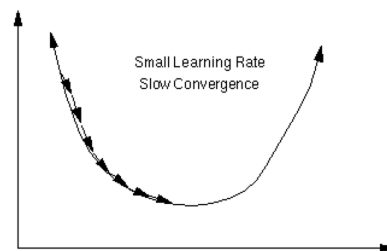
- Large learning rate

- “Overshoots”



- Small learning rate

- Extremely slow



- Solution

- Start with large  $\eta$  and settle on optimal value
- Need a schedule for shrinking  $\eta$



# Convergence of Steepest Descent

- Steepest descent converges when every element of the gradient is zero
  - In practice, very close to zero
- We may be able to avoid iterative algorithm and jump to the critical point by solving the equation  $\nabla_x f(\mathbf{x}) = 0$  for  $\mathbf{x}$

## Choosing $\eta$ : Line Search

- We can choose  $\eta$  in several different ways
- Popular approach: set  $\eta$  to a small constant
- Another approach is called *line search*:
- Evaluate  $f(\mathbf{x} - \eta \nabla_{\mathbf{x}} f(\mathbf{x}))$  for several values of  $\epsilon$  and choose the one that results in smallest objective function value

# Line Search

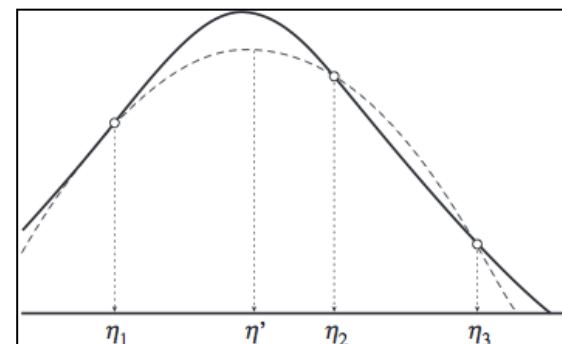
- Alternative to Simple Gradient Ascent

$$\theta^{t+1} \leftarrow \theta^t + \eta \nabla f(\theta^t)$$

- Intuition: Adaptively choose  $\eta$  at each step

- Choose direction to ascend and continue in direction until we start to descend
- Define “line” in direction of gradient

$$g(\eta) = \vec{\theta}^t + \eta \nabla f(\theta^t)$$



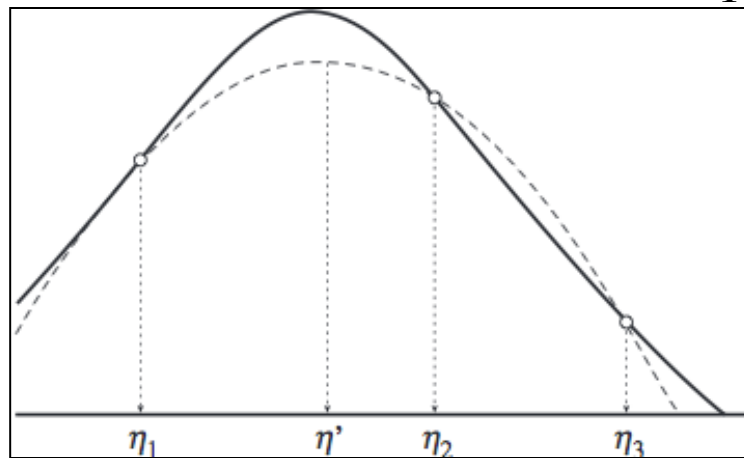
- Find three points  $\eta_1 < \eta' < \eta_3$  so that  $f(g(\eta'))$  is larger than at both  $f(g(\eta_1))$  and  $f(g(\eta_2))$

To find  $\eta'$  use binary search

Choose  $\eta' = (\eta_1 + \eta_2) / 2$

# Brent's Method

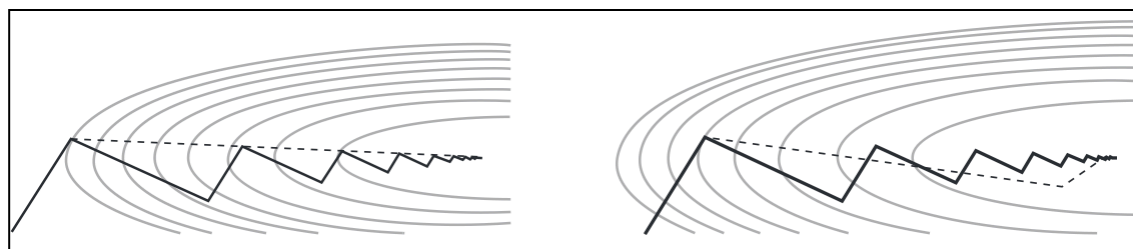
- Faster than line search
- Use quadratic function
  - Based on values of  $f$  at  $\eta_1, \eta_2, \eta_3$



- Solid line is  $f(g(\eta))$
- Three points bracket its maximum
- Dashed line shows quadratic fit

# Conjugate Gradient Descent

- In simple gradient ascent:
  - two consecutive gradients are orthogonal:
  - $\nabla f(\theta^{t+1})$  is orthogonal to  $\nabla f(\theta^t)$
  - Progress is zig-zag: progress to maximum slows down



Quadratic:  $f(x,y)=-(x^2+10y^2)$

Exponential:  $f(x,y)=\exp [-(x^2+10y^2)]$

- Solution:
  - Remember past directions and bias new direction  $\mathbf{h}^t$  as combination of current  $\mathbf{g}^t$  and past  $\mathbf{h}^{t-1}$
  - Faster convergence than simple gradient ascent

# Sequential (On-line) Learning

- Maximum likelihood solution is

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

- It is a batch technique
  - Processing entire training set in one go
- It is computationally expensive for large data sets
  - Due to huge  $N \times M$  Design matrix  $\Phi$
- Solution is to use a sequential algorithm where samples are presented one at a time
- Called Stochastic gradient descent

# Stochastic Gradient Descent

- Error function  $E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{ t_n - \mathbf{w}^T \phi(\mathbf{x}_n) \}^2$  sums over data
  - Denoting  $E_D(\mathbf{w}) = \sum_n E_n$  update parameter vector  $\mathbf{w}$  using

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n$$

- Substituting for the derivative

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \eta (t_n - \mathbf{w}^{(\tau)T} \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n)$$

- $\mathbf{w}$  is initialized to some starting vector  $\mathbf{w}^{(0)}$
- $\eta$  chosen with care to ensure convergence
- Known as *Least Mean Squares* Algorithm

# SGD Performance

Most practitioners use SGD for DL

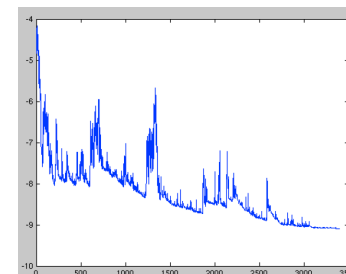
Consists of showing the input vector for a few examples,  
computing the outputs and the errors,  
Computing the *average gradient* for those examples,  
and adjusting the weights accordingly.

Process repeated for many small sets of examples  
from the training set until the average of the objective function stops decreasing.

Called stochastic because each small set of examples  
gives a noisy estimate of the average gradient over all examples.

Usually finds a good set of weights quickly compared to elaborate optimization techniques.

After training, performance is measured  
on a different test set:  
tests generalization ability of the machine  
— its ability to produce sensible answers on new inputs  
never seen during training.



Fluctuations in objective  
as gradient steps are taken  
in mini batches

24