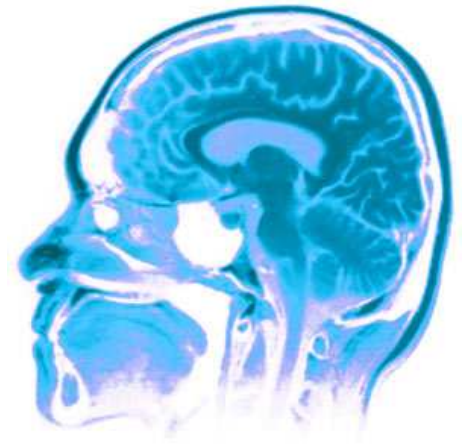




# CPSC540



## Optimization: gradient descent and Newton's method



Nando de Freitas  
*November, 2012*  
*University of British Columbia*

# Outline of the lecture

Many machine learning problems can be cast as optimization problems. This lecture introduces optimization. The objective is for you to learn:

- ☐ The definitions of gradient and Hessian.
- ☐ The gradient descent algorithm.
- ☐ Newton's algorithm.
- ☐ The stochastic gradient descent algorithm for online learning.
- ☐ How to apply all these algorithms to linear regression.

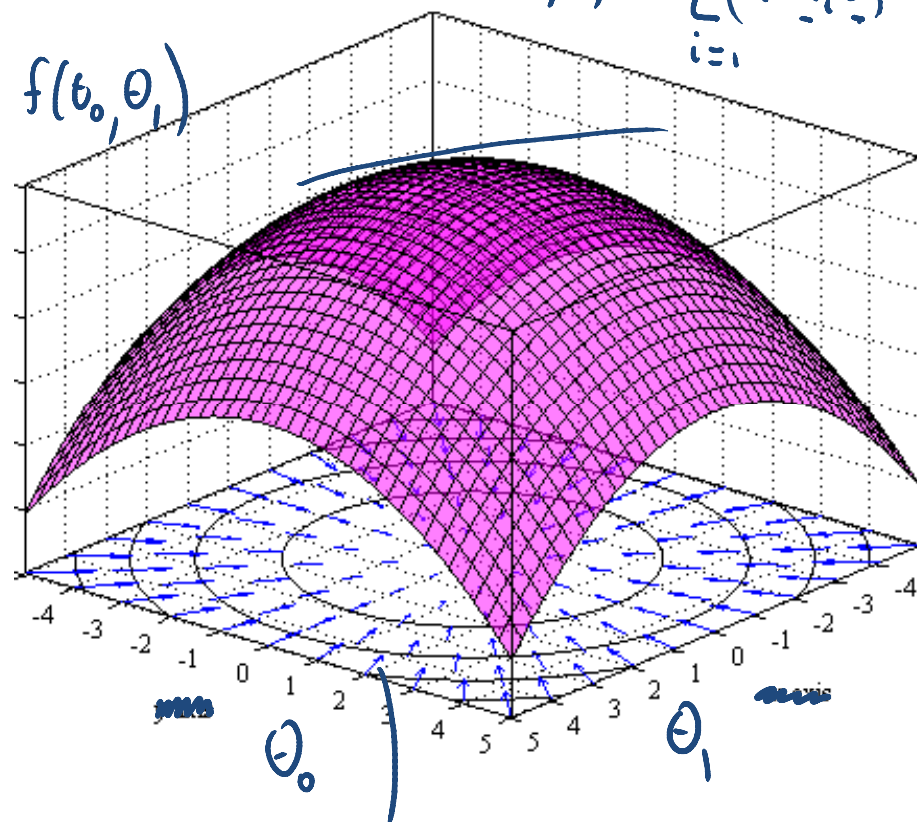
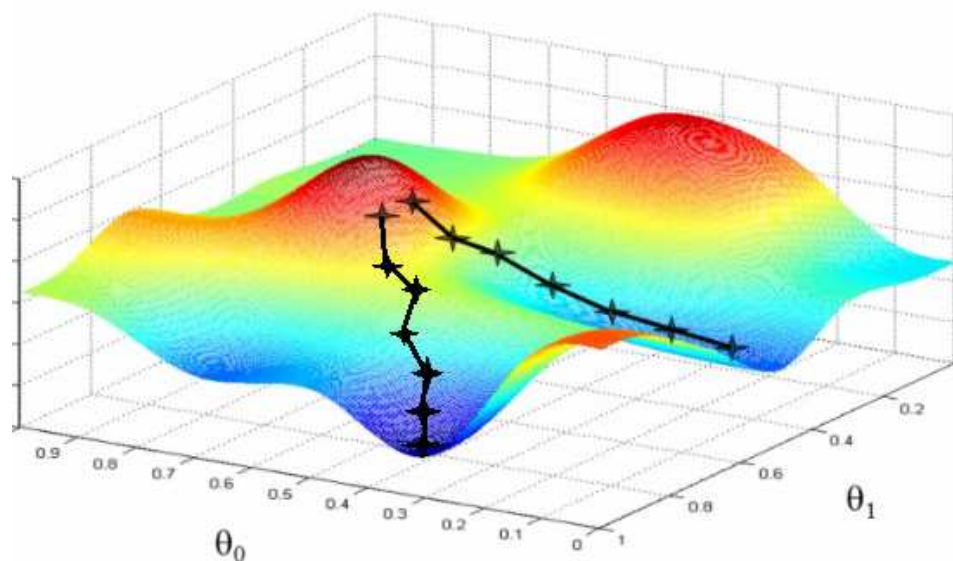
# Gradient vector

$$\boldsymbol{\theta} = (\theta_0, \theta_1)$$

Let  $\boldsymbol{\theta}$  be an  $d$ -dimensional vector and  $f(\boldsymbol{\theta})$  a scalar-valued function. The gradient vector of  $f(\cdot)$  with respect to  $\boldsymbol{\theta}$  is:

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = \begin{bmatrix} \frac{\partial f(\boldsymbol{\theta})}{\partial \theta_1} \\ \frac{\partial f(\boldsymbol{\theta})}{\partial \theta_2} \\ \vdots \\ \frac{\partial f(\boldsymbol{\theta})}{\partial \theta_n} \end{bmatrix}$$

$$f(\theta_0, \theta_1) = - \sum_{i=1}^n (y_i - x_i \theta_0)^2$$



$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = \begin{bmatrix} \frac{\partial f(\boldsymbol{\theta})}{\partial \theta_0} \\ \frac{\partial f(\boldsymbol{\theta})}{\partial \theta_1} \end{bmatrix}_{\boldsymbol{\theta} = \boldsymbol{\theta}^*}$$

# Hessian matrix

The **Hessian** matrix of  $f(\cdot)$  with respect to  $\boldsymbol{\theta}$ , written  $\nabla_{\boldsymbol{\theta}}^2 f(\boldsymbol{\theta})$  or simply as  $\mathbf{H}$ , is the  $d \times d$  matrix of partial derivatives,

$$\nabla_{\boldsymbol{\theta}}^2 f(\boldsymbol{\theta}) = \begin{bmatrix} \frac{\partial^2 f(\boldsymbol{\theta})}{\partial \theta_1^2} & \frac{\partial^2 f(\boldsymbol{\theta})}{\partial \theta_1 \partial \theta_2} & \cdots & \frac{\partial^2 f(\boldsymbol{\theta})}{\partial \theta_1 \partial \theta_d} \\ \frac{\partial^2 f(\boldsymbol{\theta})}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 f(\boldsymbol{\theta})}{\partial \theta_2^2} & \cdots & \frac{\partial^2 f(\boldsymbol{\theta})}{\partial \theta_2 \partial \theta_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\boldsymbol{\theta})}{\partial \theta_d \partial \theta_1} & \frac{\partial^2 f(\boldsymbol{\theta})}{\partial \theta_d \partial \theta_2} & \cdots & \frac{\partial^2 f(\boldsymbol{\theta})}{\partial \theta_d^2} \end{bmatrix}$$

In **offline** learning, we have a **batch** of data  $\mathbf{x}_{1:n} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ . We typically optimize cost functions of the form

$$f(\boldsymbol{\theta}) = f(\boldsymbol{\theta}, \mathbf{x}_{1:n}) = \frac{1}{n} \sum_{i=1}^n f(\boldsymbol{\theta}, \mathbf{x}_i)$$

The corresponding gradient is

$$g(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}, \mathbf{x}_i)$$

For linear regression with training data  $\{\mathbf{x}_i, y_i\}_{i=1}^n$ , we have have the quadratic cost

$$f(\boldsymbol{\theta}) = f(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) = \sum_{i=1}^n (y_i - \mathbf{x}_i \boldsymbol{\theta})^2$$

# Gradient vector and Hessian matrix

$$f(\boldsymbol{\theta}) = f(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y}) = \underbrace{(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})}_{\underbrace{\sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\theta})^2}}$$

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = 2\mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - 2\mathbf{X}^T \mathbf{y} \equiv \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = -2 \sum_{i=1}^n \mathbf{x}_i^T (y_i - \mathbf{x}_i^T \boldsymbol{\theta})$$

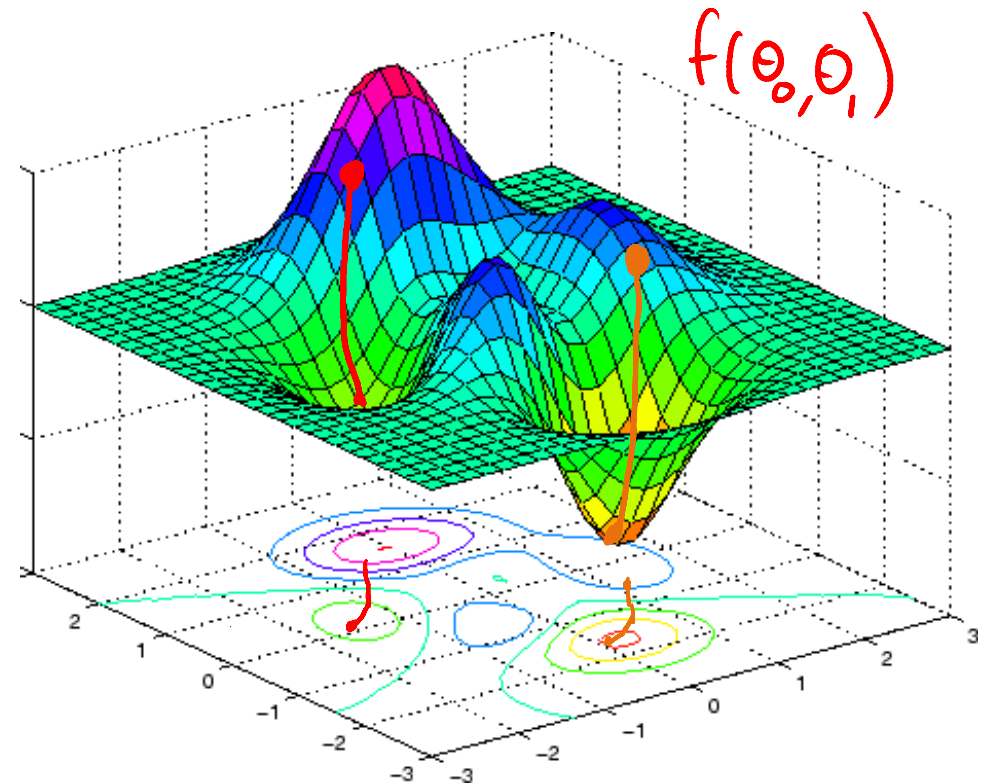
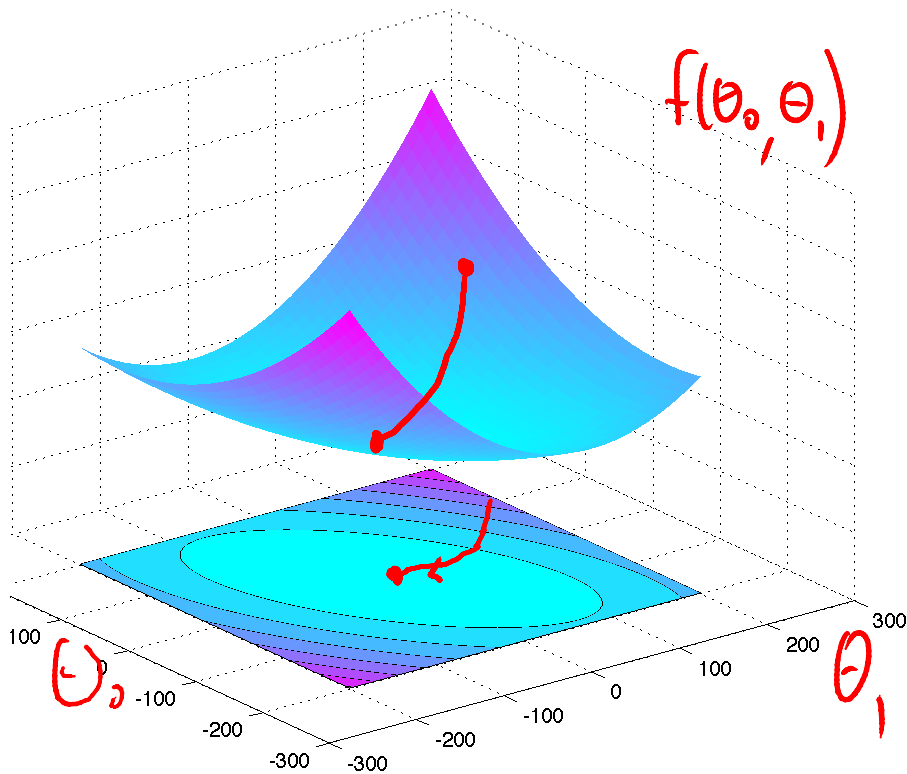
$$H = \nabla_{\boldsymbol{\theta}}^2 f(\boldsymbol{\theta}) = 2\mathbf{X}^T \mathbf{X}$$

# Steepest gradient descent algorithm

One of the simplest optimization algorithms is called **gradient descent** or **steepest descent**. This can be written as follows:

$$\theta_{k+1} = \theta_k - \eta_k \mathbf{g}_k = \theta_k - \eta_k \nabla f(\theta_k)$$

where  $k$  indexes steps of the algorithm,  $\mathbf{g}_k = \mathbf{g}(\theta_k)$  is the gradient at step  $k$ , and  $\eta_k > 0$  is called the **learning rate** or **step size**.



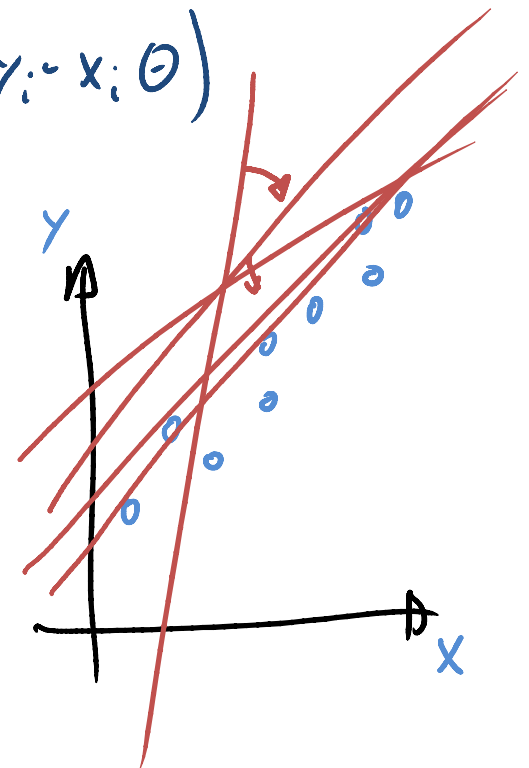
# Steepest gradient descent algorithm for least squares

$$f(\boldsymbol{\theta}) = f(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) = \sum_{i=1}^n (y_i - \mathbf{x}_i \boldsymbol{\theta})^2$$

$$\nabla f(\boldsymbol{\theta}) = 2 \mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - 2 \mathbf{X}^T \mathbf{y} = -2 \sum_{i=1}^n \mathbf{x}_i^T (y_i - \mathbf{x}_i \boldsymbol{\theta})$$

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - 2\eta \left[ \mathbf{X}^T \mathbf{X} \boldsymbol{\theta}_k - \mathbf{X}^T \mathbf{y} \right]$$

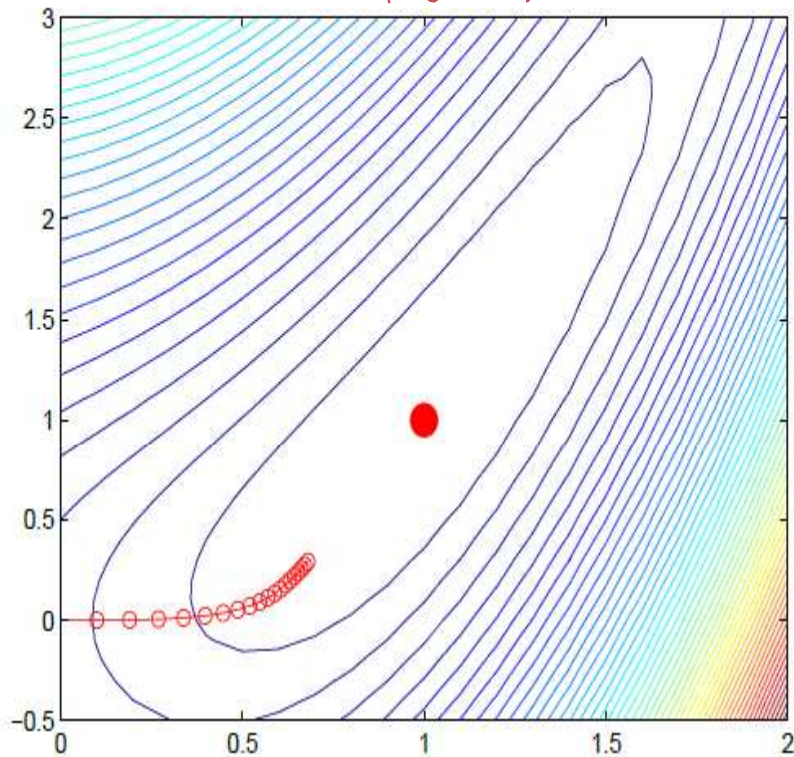
$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + 2\eta \sum_{i=1}^n \mathbf{x}_i^T (y_i - \mathbf{x}_i \boldsymbol{\theta})$$





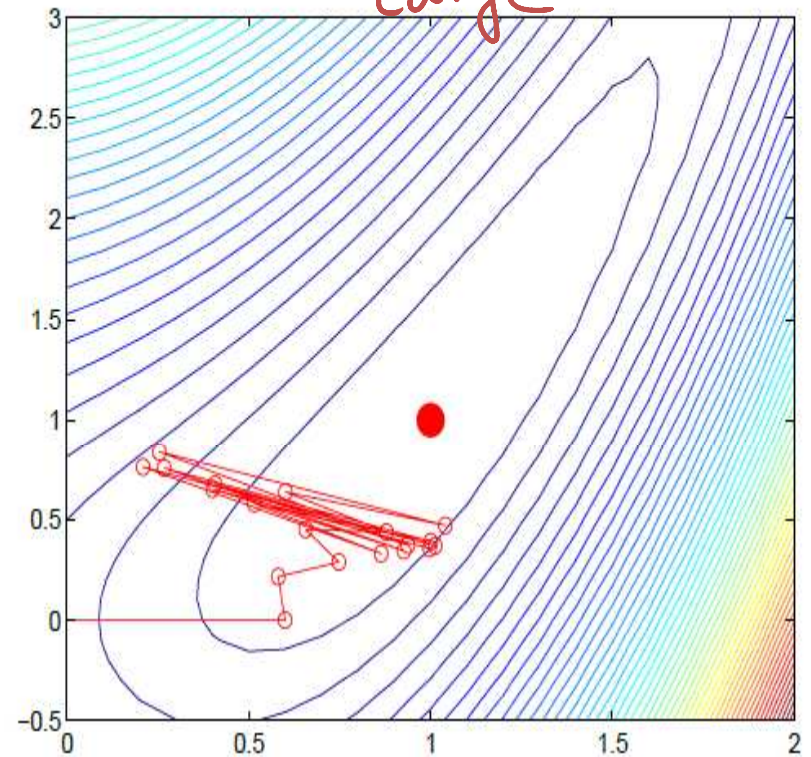
# How to choose the step size ?

Too small



$\eta = 0.1.$

Too large



$\eta = 0.6.$

$$\Theta_{k+1} = \Theta_k - \underbrace{(0.1) \nabla f(\Theta_k)}_{\approx 0}$$

# Newton's algorithm

The most basic second-order optimization algorithm is **Newton's algorithm**, which consists of updates of the form

$$\theta_{k+1} = \theta_k - \mathbf{H}_k^{-1} \mathbf{g}_k$$

This algorithm is derived by making a second-order Taylor series approximation of  $f(\theta)$  around  $\theta_k$ :

$$\underline{f_{quad}(\theta)} = \underline{f(\theta_k)} + \underline{\mathbf{g}_k^T (\theta - \theta_k)} + \frac{1}{2} (\theta - \theta_k)^T \mathbf{H}_k (\theta - \theta_k)$$

$\frac{1}{2} [\theta^T \mathbf{H}_k \theta - 2\theta^T \mathbf{H}_k \theta_k + \cancel{\theta_k^T \mathbf{H}_k \theta_k}]$

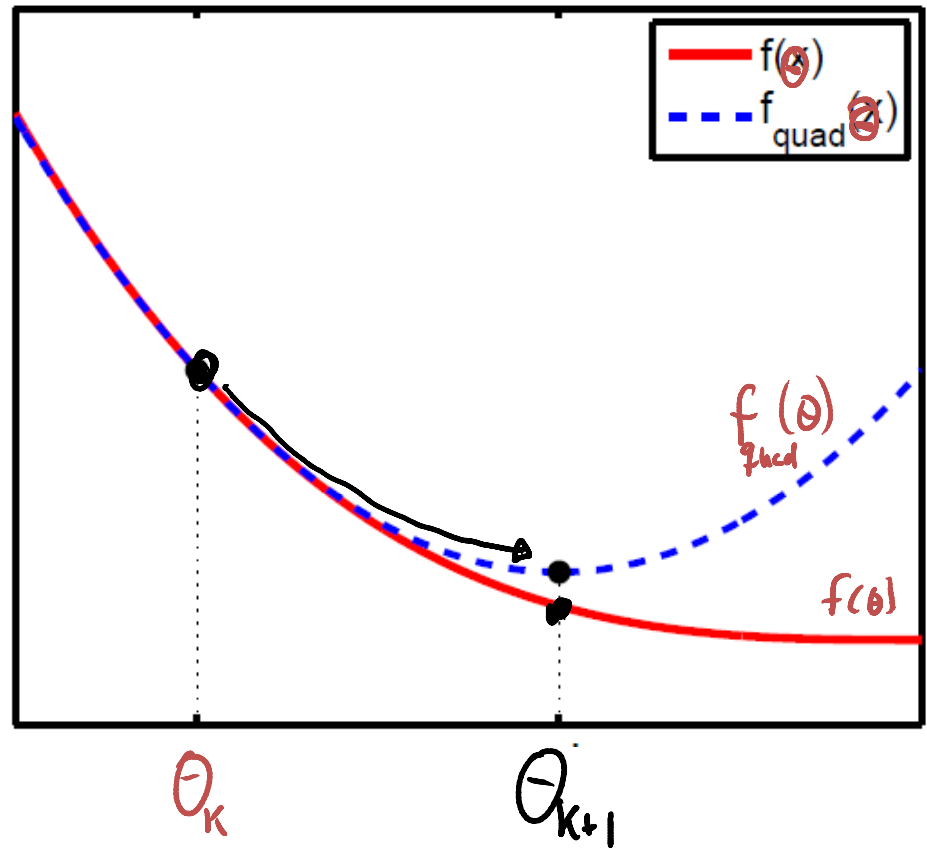
differentiating and equating to zero to solve for  $\theta_{k+1}$ .

$$\frac{\partial f_{quad}(\theta)}{\partial \theta} = 0 + \mathbf{g}_k + \mathbf{H}_k (\theta - \theta_k) = 0$$

$$-\mathbf{g}_k = \mathbf{H}_k \theta - \mathbf{H}_k \theta_k$$

$$-\mathbf{H}_k^{-1} \mathbf{g}_k = \theta - \theta_k$$

# Newton's as bound optimization



# Newton's algorithm for linear regression

$$f(\boldsymbol{\theta}) = f(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) = \sum_{i=1}^n (y_i - \mathbf{x}_i \boldsymbol{\theta})^2$$

$$\mathbf{g} = \nabla f(\boldsymbol{\theta}) = 2\mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - 2\mathbf{X}^T \mathbf{y}$$

$$\mathbf{H} = 2\mathbf{X}^T \mathbf{X}$$

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \frac{1}{2} (\mathbf{X}^T \mathbf{X})^{-1} [\mathbf{X}^T \mathbf{X} \boldsymbol{\theta}_k - \mathbf{X}^T \mathbf{y}]$$

$$= \boldsymbol{\theta}_k - (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}) \boldsymbol{\theta}_k + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Exercise  
Show This is  
 $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

# Advanced: Newton CG algorithm

Rather than computing  $\mathbf{d}_k = -\mathbf{H}_k^{-1} \mathbf{g}_k$  directly, we can solve the linear system of equations  $\mathbf{H}_k \mathbf{d}_k = -\mathbf{g}_k$  for  $\mathbf{d}_k$ .

One efficient and popular way to do this, especially if  $\mathbf{H}$  is sparse, is to use a conjugate gradient method to solve the linear system.

---

```
1 Initialize  $\boldsymbol{\theta}_0$ 
2 for  $k = 1, 2, \dots$  until convergence do
3   Evaluate  $\mathbf{g}_k = \nabla f(\boldsymbol{\theta}_k)$ 
4   Evaluate  $\mathbf{H}_k = \nabla^2 f(\boldsymbol{\theta}_k)$ 
5   Solve  $\mathbf{H}_k \mathbf{d}_k = -\mathbf{g}_k$  for  $\mathbf{d}_k$ 
6   Use line search to find stepsize  $\eta_k$  along  $\mathbf{d}_k$ 
7    $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \eta_k \mathbf{d}_k$ 
```

---

# Online learning aka stochastic gradient descent

Batch

$$\Theta_{k+1} = \Theta_k + \eta \sum_{i=1}^n x_i^T (y_i - x_i \Theta_k) \quad \left( \begin{smallmatrix} n \\ \text{data} \\ \text{points} \end{smallmatrix} \right)$$

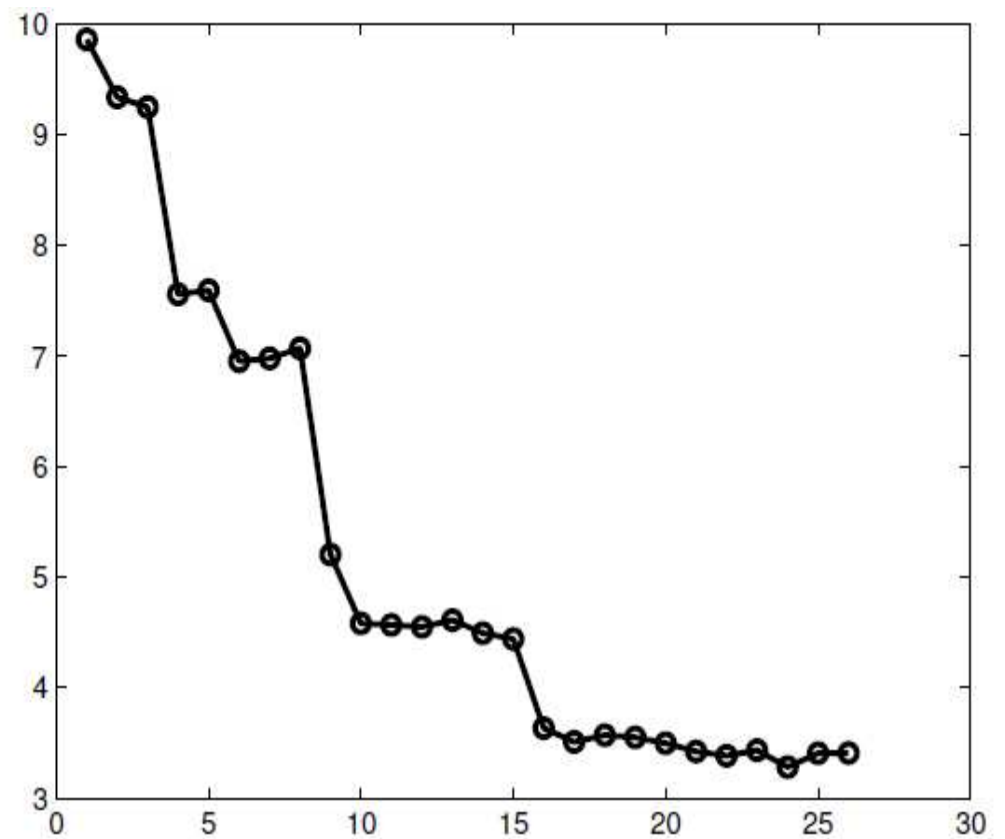
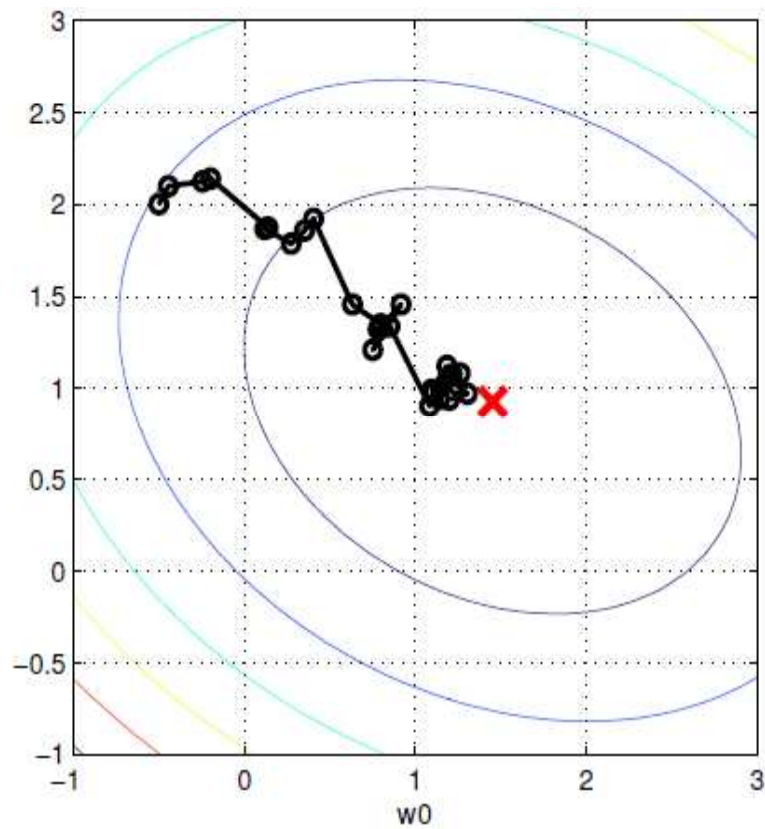
Online

$$\Theta_{k+1} = \Theta_k + \eta x_k^T (y_k - x_k \Theta_k)$$

mini-batch

$$\Theta_{k+1} = \Theta_k + \eta \sum_{j=1}^{20} x_j^T (y_j - x_j \Theta_k)$$

# The online learning algorithm





# Stochastic gradient descent

SGD can also be used for offline learning, by repeatedly cycling through the data; each such pass over the whole dataset is called an **epoch**. This is useful if we have **massive datasets** that will not fit in main memory. In this offline case, it is often better to compute the gradient of a **mini-batch** of  $B$  data cases. If  $B = 1$ , this is standard SGD, and if  $B = N$ , this is standard steepest descent. Typically  $B \sim 100$  is used.

Intuitively, one can get a fairly good estimate of the gradient by looking at just a few examples. Carefully evaluating precise gradients using large datasets is often a waste of time, since the algorithm will have to recompute the gradient again anyway at the next step. It is often a better use of computer time to have a noisy estimate and to move rapidly through parameter space.

SGD is often less prone to getting stuck in shallow local minima, because it adds a certain amount of “noise”. Consequently it is quite popular in the machine learning community for fitting models such as neural networks and deep belief networks with non-convex objectives.



# Next lecture

In the next lecture, we apply these ideas to learn a neural network with a single neuron (logistic regression).