

Outline of the lecture

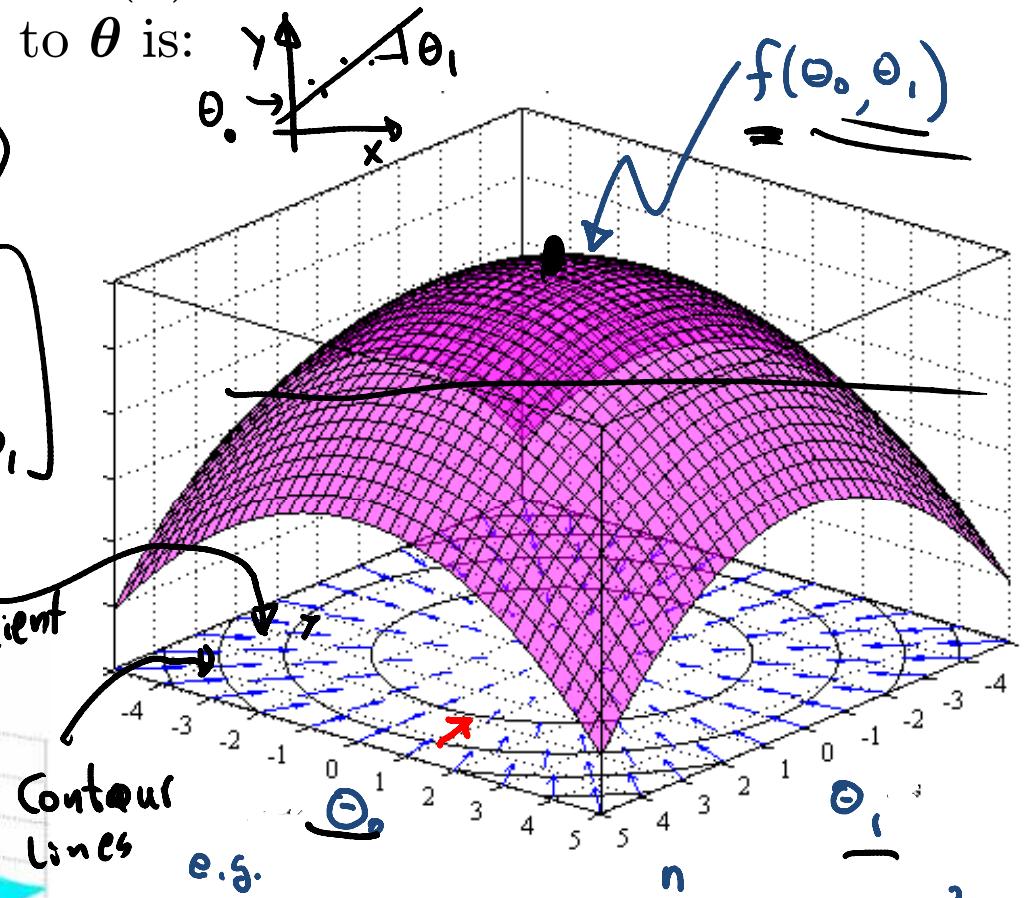
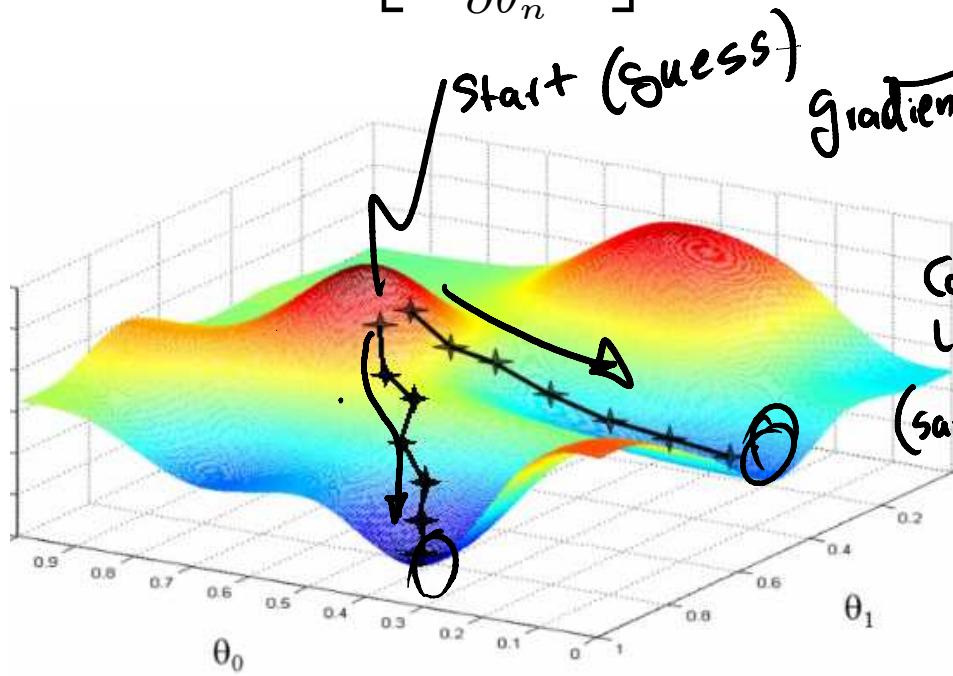
Many machine learning problems can be cast as optimization problems. This lecture introduces optimization. The objective is for you to learn:

- 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 ∇f

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

$$\nabla_{\theta} f(\theta) = \begin{bmatrix} \frac{\partial f(\theta)}{\partial \theta_1} \\ \frac{\partial f(\theta)}{\partial \theta_2} \\ \vdots \\ \frac{\partial f(\theta)}{\partial \theta_n} \end{bmatrix} \quad \approx \quad \begin{bmatrix} \nabla_{\theta} f(\theta_0, \theta_1) \\ \frac{\partial f}{\partial \theta_0} \\ \frac{\partial f}{\partial \theta_1} \end{bmatrix}$$



$$(same \ height) \quad f(\theta_0, \theta_1) = - \sum_{i=1}^n (\gamma_i - \underline{x}_i \cdot \underline{\theta})^2$$

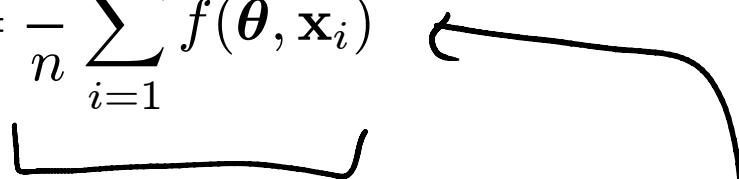
$$\underline{\theta} = (\theta_0, \theta_1) \in \mathbb{R}^2$$

Hessian matrix

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

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

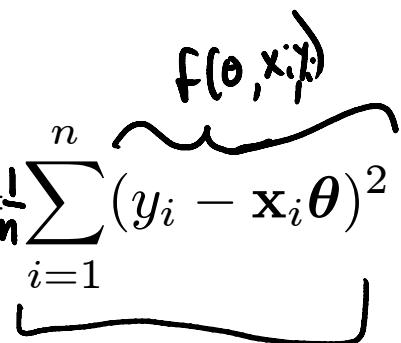
In **offline** learning, we have a **batch** of data $\mathbf{x}_{1:n} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$. We typically optimize cost functions of the form

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


The corresponding gradient is

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

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

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


Gradient vector and Hessian matrix

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

$$\begin{aligned}\nabla f(\theta) &= \frac{\partial}{\partial \theta} (\mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\theta + \theta^\top \mathbf{X}^\top \mathbf{X}\theta) \\ &= -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\theta \quad \stackrel{=} \Bigg| \text{DF}(\theta) = -2 \sum_{i=1}^n \mathbf{x}_i^\top (y_i - \mathbf{x}_i^\top \theta)\end{aligned}$$

$$\nabla^2 f(\theta) = 0 + 2\mathbf{X}^\top \mathbf{X}$$

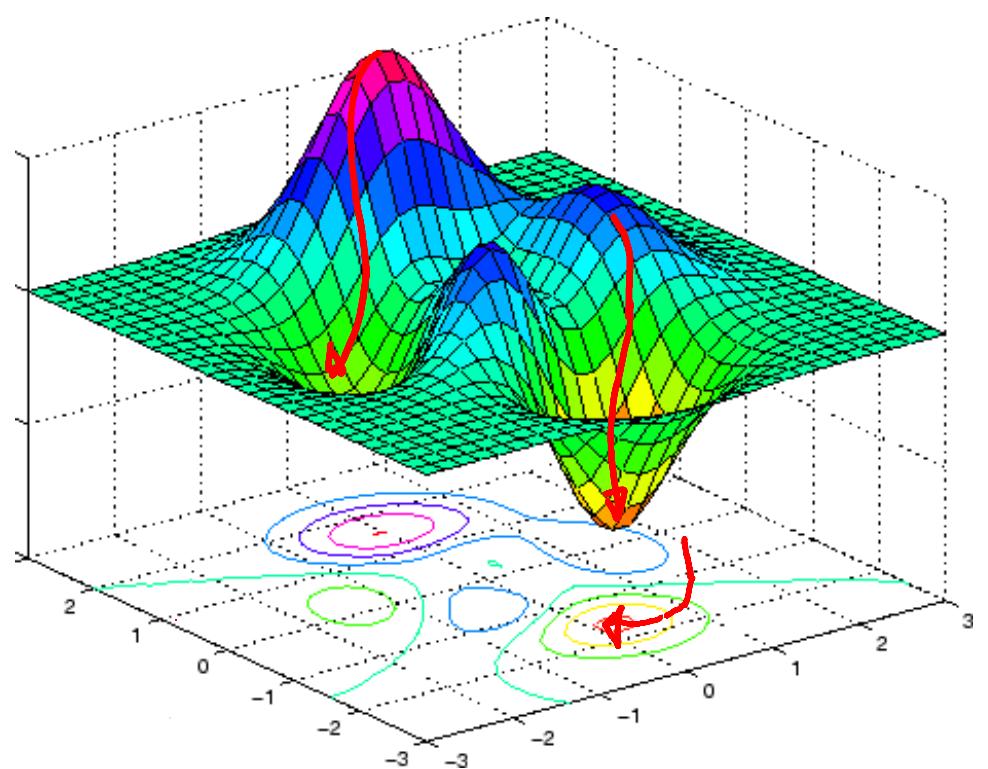
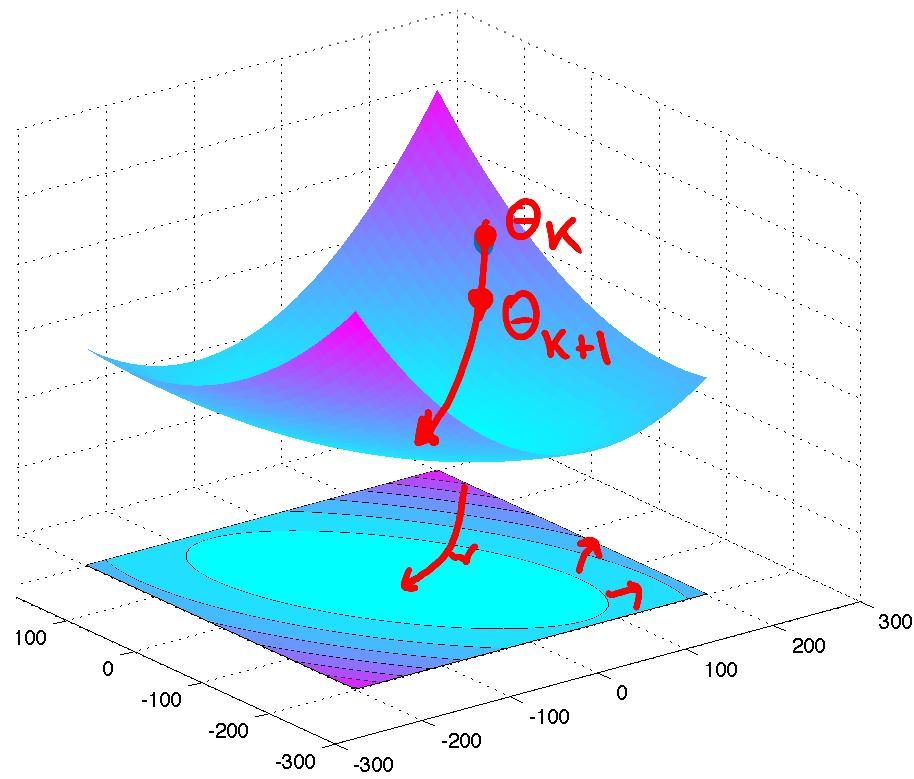
$$= 2\mathbf{X}^\top \mathbf{X}$$

Steepest gradient descent algorithm

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

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

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

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

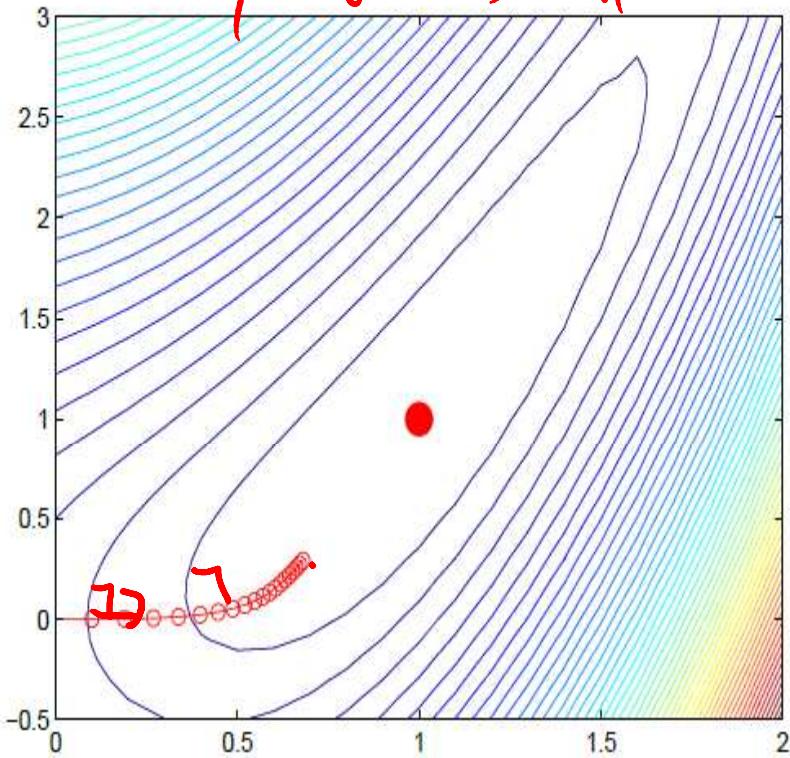
$$\nabla f(\theta) = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \theta$$

$$\theta_{k+1} = \theta_k - \eta \left[\underbrace{-2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \theta_k}_{\text{gradient}} \right]$$

$$\theta_{k+1} = \theta_k - \eta \left[-2 \sum_{i=1}^n \mathbf{x}_i^\top (\mathbf{y}_i - \mathbf{x}_i^\top \theta_k) \right]$$

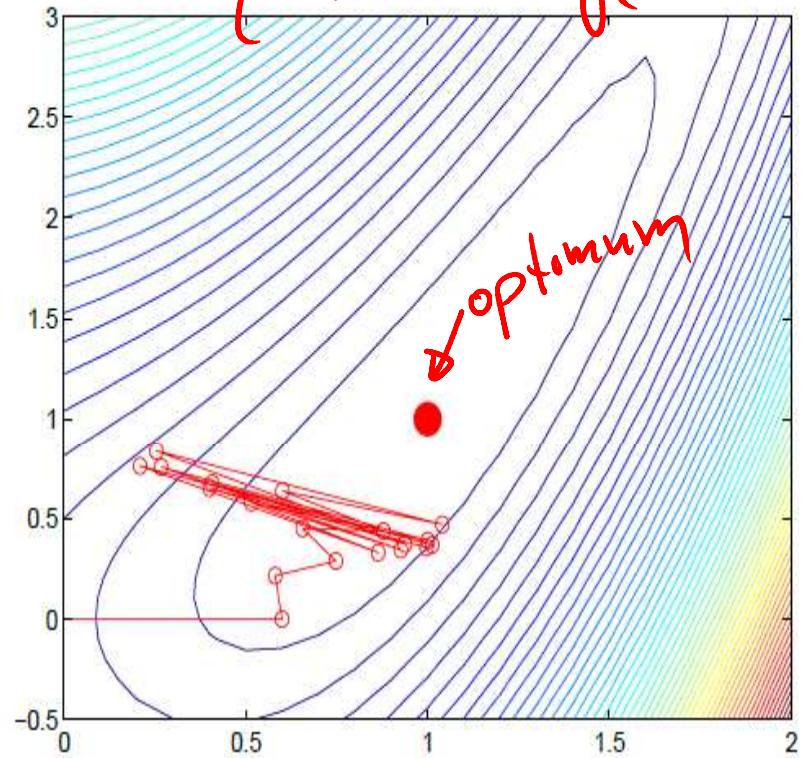
How to choose the step size ?

η too small



$$\eta = \underline{0.1}$$

η too large



$$\eta = 0.6$$

$$\Theta_{K+1} = \Theta_K - \eta \nabla f(\Theta_K)$$

Newton's algorithm

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

$$\theta_{k+1} = \theta_k - H_K^{-1} g_k$$

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

$$f_{quad}(\theta) = f(\underline{\theta_k}) + \underbrace{g_k^T(\theta - \theta_k)}_{\text{new value}} + \frac{1}{2} (\theta - \theta_k)^T H_k (\theta - \theta_k)$$

old value $\Delta\theta$

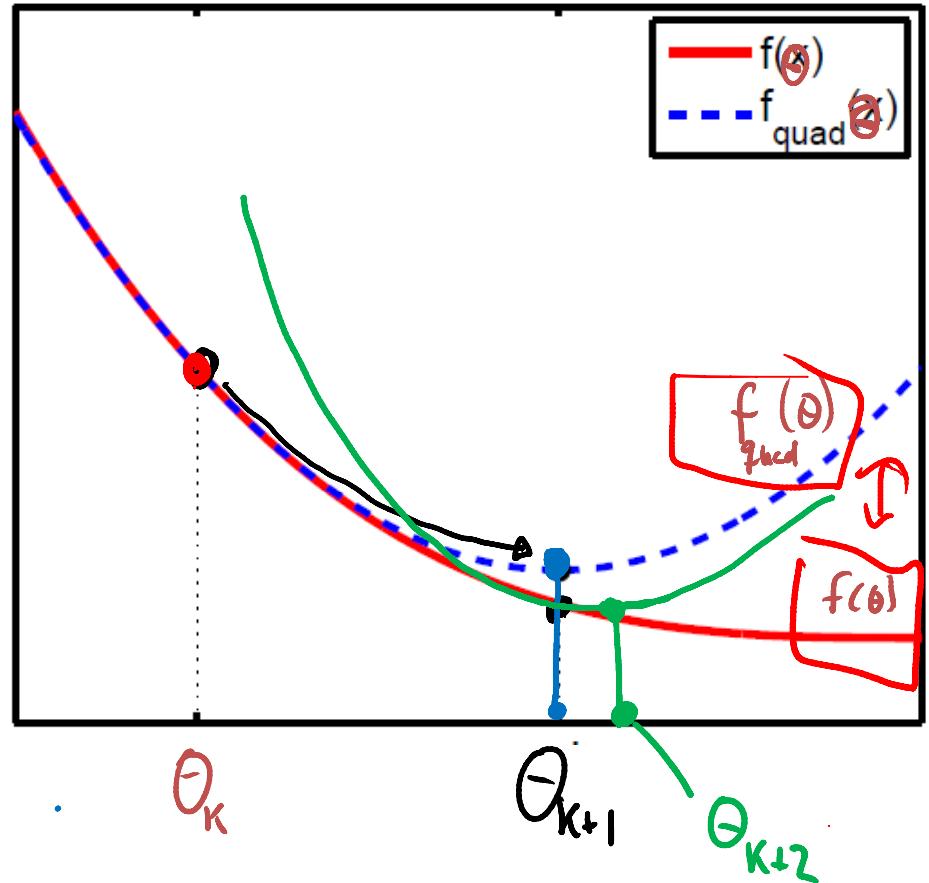
differentiating and equating to zero to solve for θ_{k+1} .

$$\nabla f_{quad}(\theta) = 0 + g_k + H_k(\theta - \theta_k) = 0$$

$$-g_k = H_k(\theta - \theta_k)$$

$$\theta = \theta_k - H_k^{-1} g_k$$

Newton's as bound optimization



Newton's algorithm for linear regression

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

$$\mathcal{S} = \nabla f(\theta) = -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\theta$$

$$\mathbf{H} = \nabla^2 f(\theta) = 2\mathbf{X}^\top \mathbf{X}$$

$$\theta_{k+1} = \theta_k - H_k^{-1} \mathcal{S}_k$$

$$\begin{aligned} &= \theta_k - [2\mathbf{X}^\top \mathbf{X}]^{-1} [-2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\theta_k] \\ &= \theta_k + \underbrace{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}}_{\text{Curly brace}} - \cancel{(\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{X})} \theta_k \end{aligned}$$

Advanced: Newton CG algorithm

$$\theta_{k+1} = \theta_k + d_k$$

Rather than computing $d_k = -H_k^{-1}g_k$ directly, we can solve the linear system of equations $H_k d_k = -g_k$ for d_k .

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

-
- 1 Initialize θ_0
 - 2 **for** $k = 1, 2, \dots$ until convergence **do**
 - 3 Evaluate $g_k = \nabla f(\theta_k)$
 - 4 Evaluate $H_k = \nabla^2 f(\theta_k)$
 - 5 Solve $H_k d_k = -g_k$ for d_k minres
 - 6 Use line search to find stepsize η_k along d_k
 - 7 $\theta_{k+1} = \theta_k + \eta_k d_k$
-

conjugate
gradient

402

Estimating the mean recursively

average = $\Theta_N = \frac{1}{N} \sum_{i=1}^N x_i$ BATCH

$$\begin{aligned}\Theta_N &= \frac{1}{N} X_N + \frac{1}{N} \frac{N-1}{N-1} \sum_{i=1}^{N-1} x_i \\ &= \frac{1}{N} X_N + \frac{1}{N-1} \left(\frac{N-1}{N} \right) \sum_{i=1}^{N-1} x_i = \frac{1}{N} X_N + \left(\frac{N-1}{N} \right) \Theta_{N-1}\end{aligned}$$

$$\Theta_N = \left(1 - \frac{1}{N} \right) \Theta_{N-1} + \frac{1}{N} X_N$$
 ONLINE

Online learning

$x^{(i)} \sim P(x)$

aka stochastic gradient descent

$$J(\theta) = \underbrace{\int J(\theta, x) P(x) dx}_{\text{expected cost}} \approx \frac{1}{N} \sum_{i=1}^N J(\theta, x_i)$$

$\underbrace{\qquad\qquad\qquad}_{\text{empirical cost risk}}$

$$\underline{DJ(\theta)} = \int \underline{DJ(\theta, x)} P(x) dx$$

$$\begin{aligned} \theta_{k+1} &= \theta_k - \eta \frac{1}{N} \sum_{i=1}^N \underline{J(\theta_k, x_i)} \stackrel{\text{Let } N=1}{=} \theta_k - \eta J(\theta_k, x_k) \\ &\stackrel{\text{true grad}}{=} \theta_k - \eta J(\theta_k) + \eta \left[J(\theta_k) - J(\theta_k, x_k) \right] / N \end{aligned}$$

$\overbrace{\qquad\qquad\qquad}^{\text{noise}}$

Online learning aka stochastic gradient descent

Batch

$$\theta_{k+1} = \theta_k + \eta \sum_{i=1}^n x_i^T (y_i - x_i \cdot \theta_k)$$

(n data points)

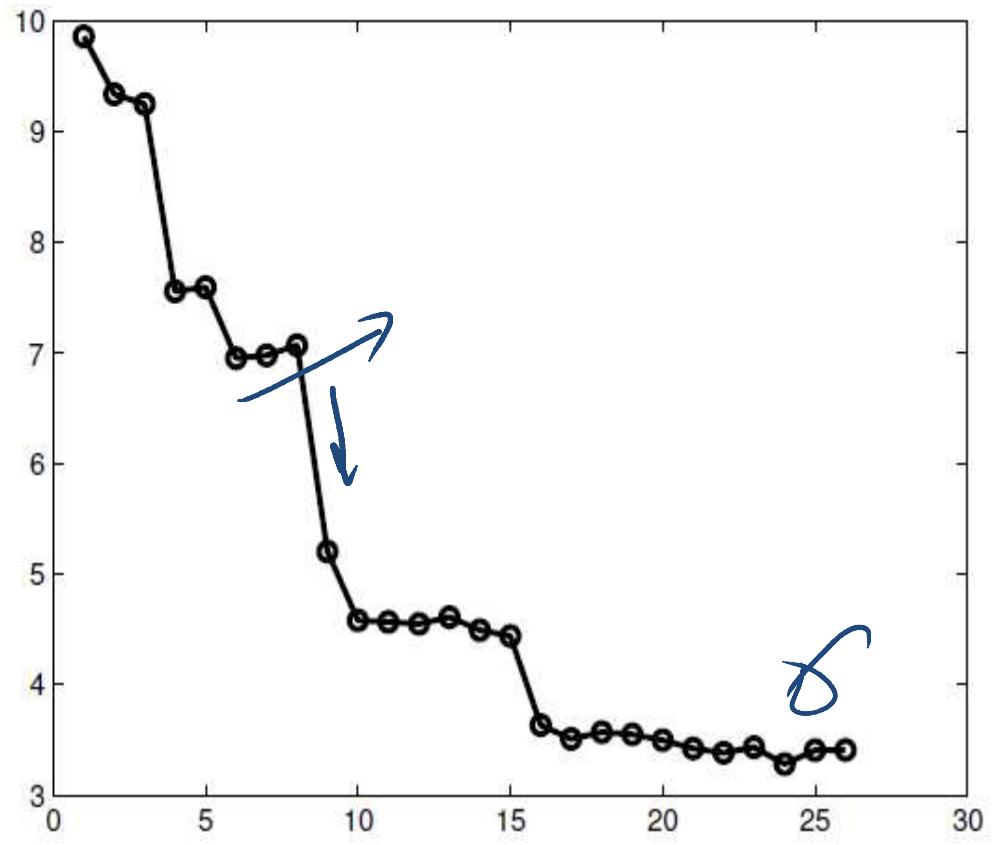
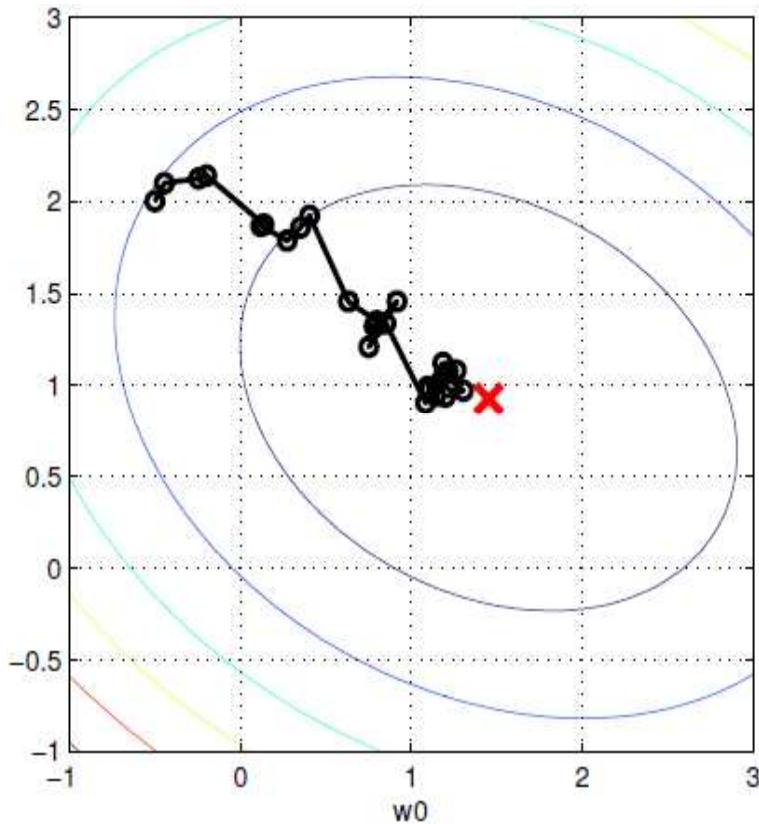
Online

$$\theta_{k+1} = \theta_k + \eta x_k^T (y_k - x_k \cdot \theta_k)$$

mini-batch

$$\theta_{k+1} = \theta_k + \eta \sum_{j=1}^{20} x_j^T (y_j - x_j \cdot \theta_k)$$

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.