



# Optimization

Nando de Freitas



UNIVERSITY OF  
**OXFORD**

# Outline of the lecture

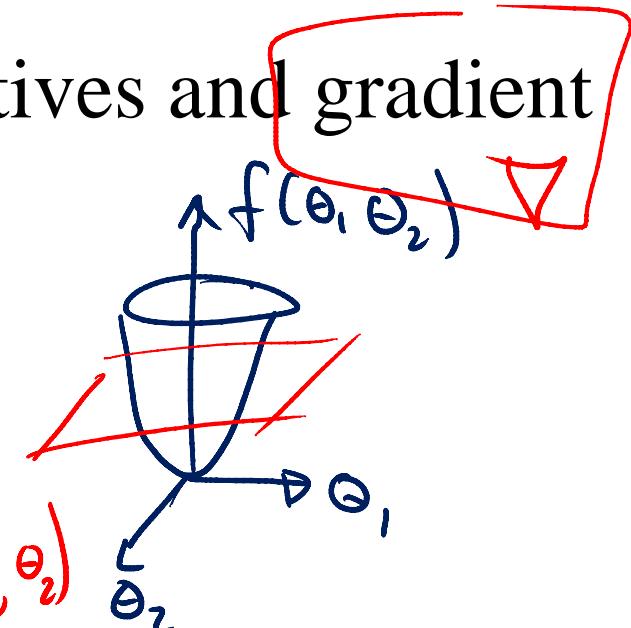
Many machine learning problems can be cast as optimization problems.  
This lecture introduces optimization. The objective is for you to learn:

- The definitions of gradient and Hessian.
- The gradient descent algorithm.
- Newton's algorithm.
- Stochastic gradient descent (SGD) for online learning.
- Popular variants, such as AdaGrad and Asynchronous SGD.
- Improvements such as momentum and Polyak averaging.
- How to apply all these algorithms to linear regression.

# Calculus background: Partial derivatives and gradient

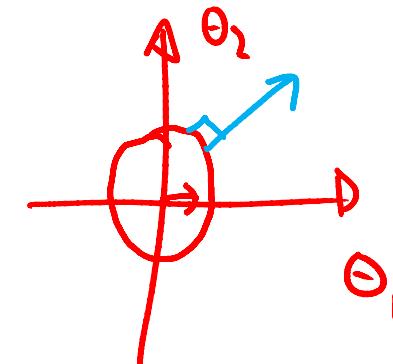
$$f(\theta) = f(\theta_1, \theta_2) = \theta_1^2 + \theta_2^2$$

$$\frac{\partial f(\theta_1, \theta_2)}{\partial \theta_1} = \lim_{\Delta \theta_1 \rightarrow 0} \frac{f(\theta_1 + \Delta \theta_1, \theta_2) - f(\theta_1, \theta_2)}{\Delta \theta_1}$$



$$\frac{\partial f(\theta)}{\partial \theta_1} = \boxed{2\theta_1}$$

$$\frac{\partial f(\theta)}{\partial \theta_2} = 2\theta_2$$



$$\nabla J(\theta) = \begin{bmatrix} 2\theta_1 \\ 2\theta_2 \end{bmatrix}$$

$$(J_1 = 1, J_2 = 1)$$

$$\nabla J(\theta_1, \theta_2) = (2, 2)$$

## Necessary calculus background: Hessian

$$\frac{\partial}{\partial \theta_1} \left( \frac{\partial f(\theta)}{\partial \theta_1} \right) = \frac{\partial^2 f(\theta)}{\partial \theta_1^2} = 2$$



$$H = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

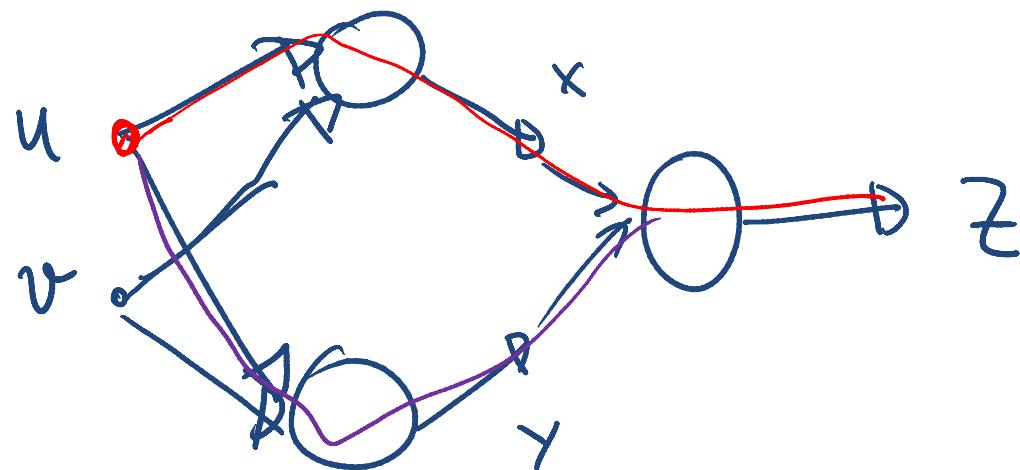
$$\frac{\partial^2 f(\theta)}{\partial \theta_2^2} = 2$$



$$\frac{\partial^2 f(\theta)}{\partial \theta_1 \partial \theta_2} = \frac{\partial^2 f(\theta)}{\partial \theta_2 \partial \theta_1} = 0$$

## Necessary calculus background: Chain rule

$$z = f(x(u,v), y(u,v))$$



$$\frac{\partial z}{\partial u} = \frac{\partial z}{\partial x} \frac{\partial x}{\partial u} + \frac{\partial z}{\partial y} \frac{\partial y}{\partial u}$$

# Necessary calculus background: Linear regression

$$y_i = \theta_0 + \theta_1 x_i \quad i=1, 2, \dots, n$$

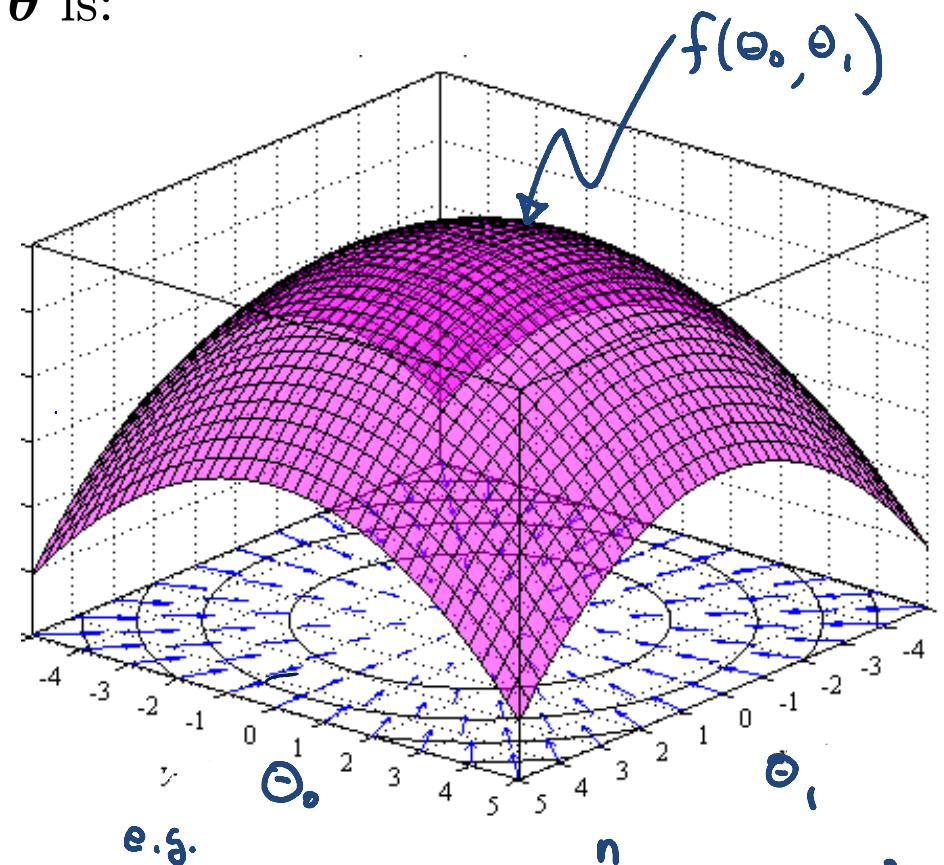
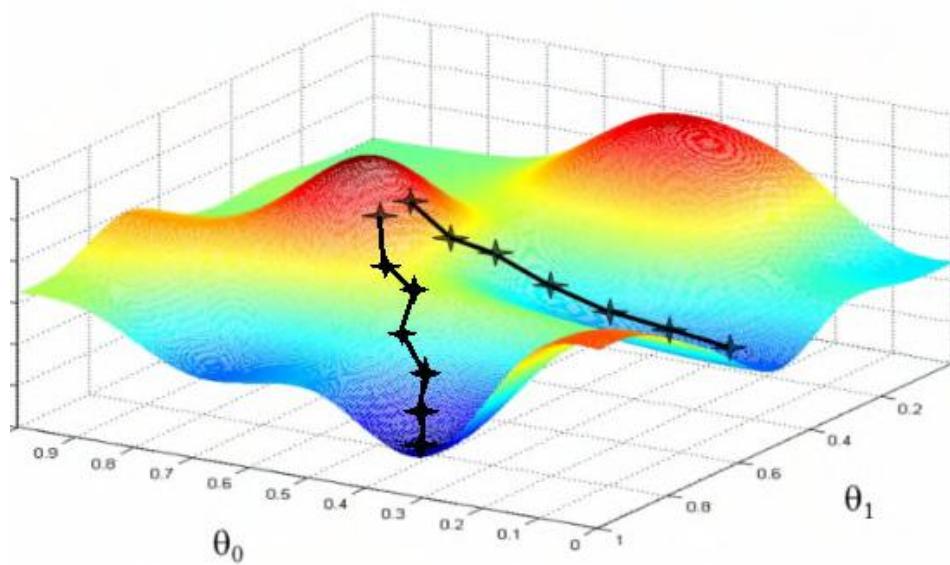
$$J(\underline{\theta}) = \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2 + \delta^2 \theta_1^2$$

$$\nabla J = \begin{bmatrix} \sum_{i=1}^n 2(y_i - \theta_0 - \theta_1 x_i)(-1) \\ \sum_{i=1}^n 2(y_i - \theta_0 - \theta_1 x_i)(x_i) + 2\delta^2 \theta_1 \end{bmatrix}$$

# Gradient vector

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

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



e.g.

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

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

# 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,

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

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

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

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 the quadratic cost

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

# Gradient vector and Hessian matrix

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

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

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

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

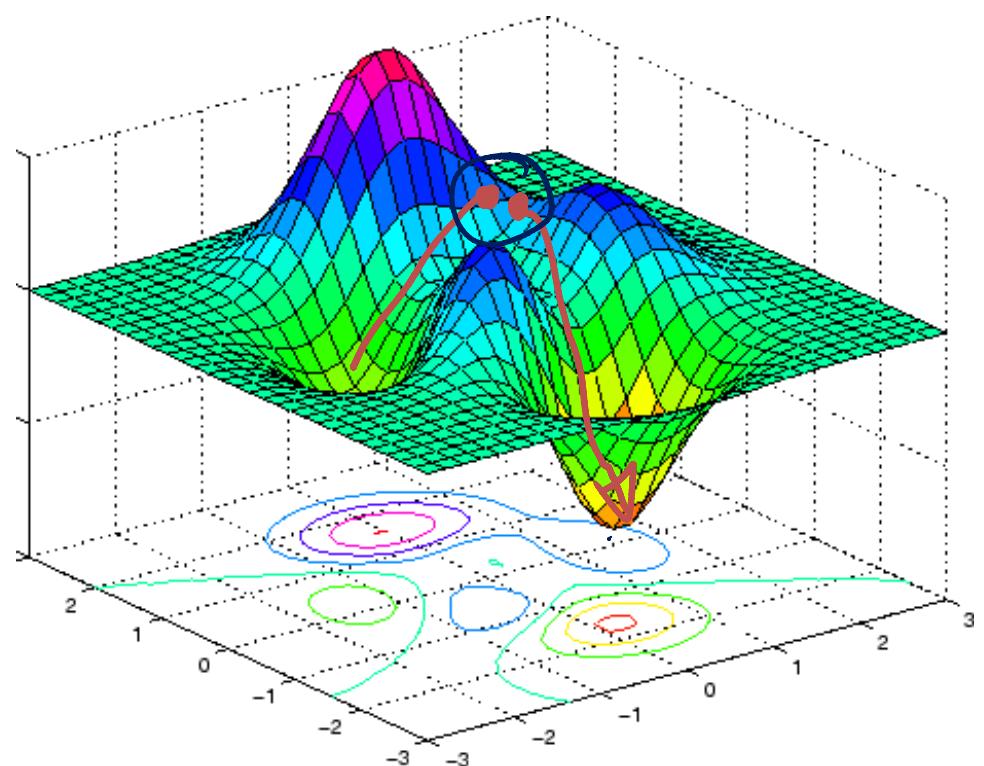
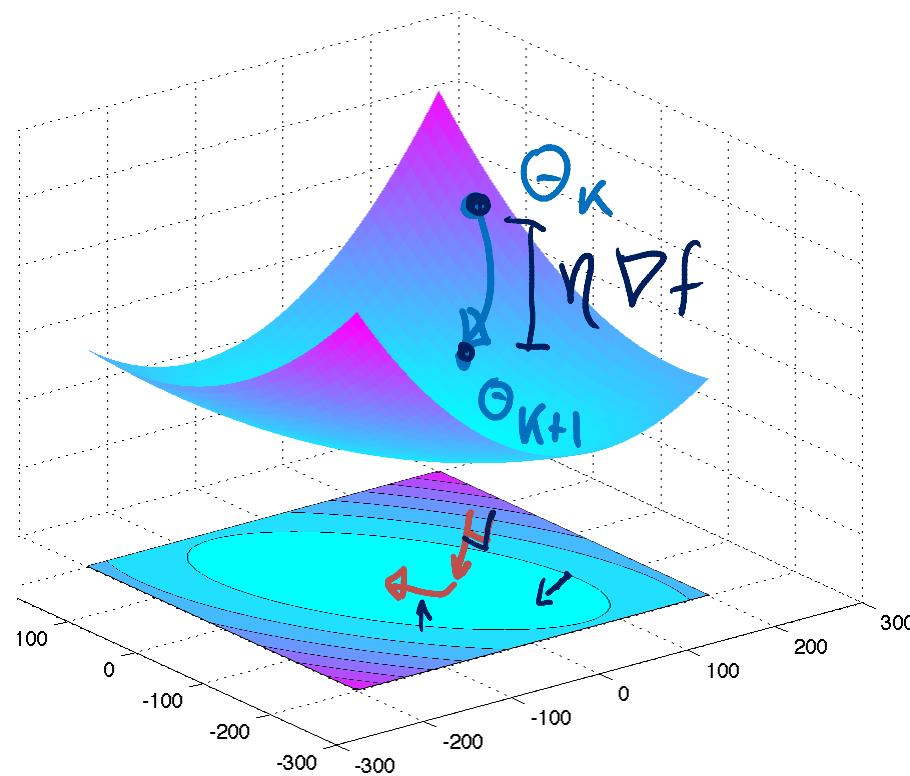
$\mathbf{X} \in \mathbb{R}^{n \times d}$

# Steepest gradient descent algorithm

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

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

where  $k$  indexes steps of the algorithm,  $g_k = g(\theta_k)$  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

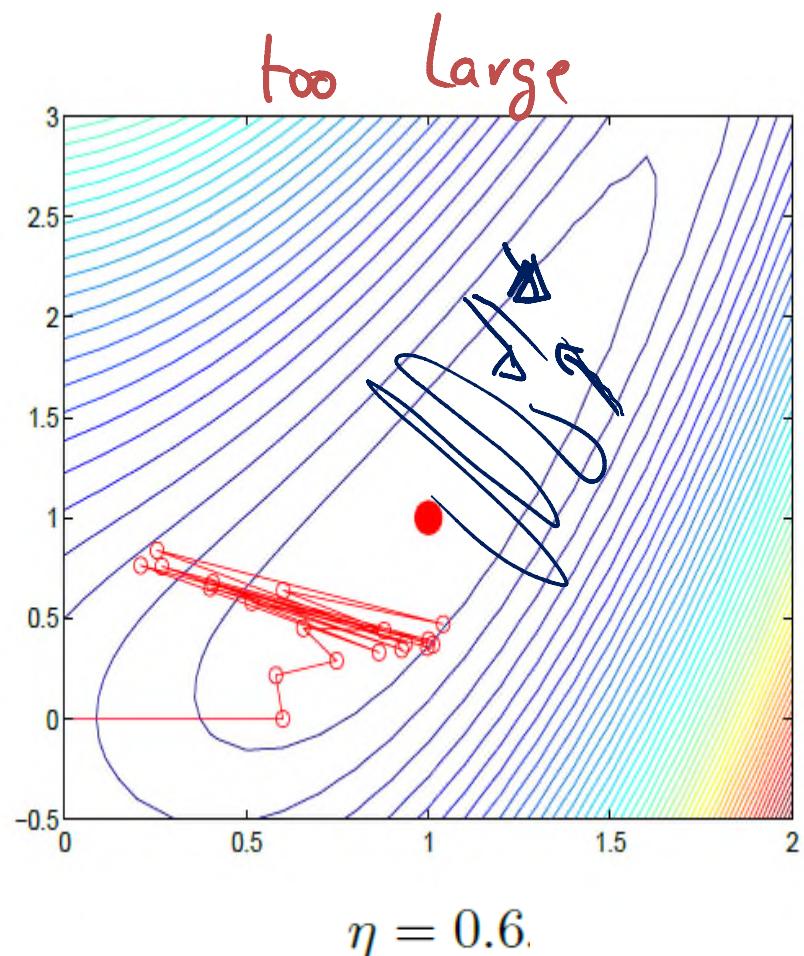
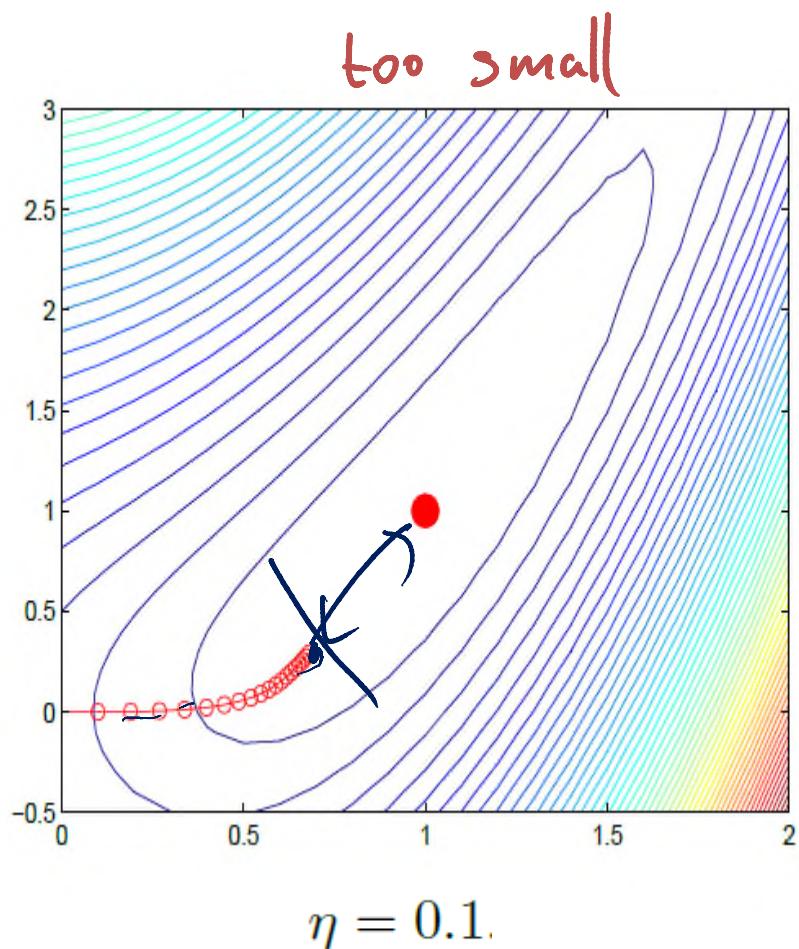
$$f(\theta) = f(\theta, \mathbf{X}, \mathbf{y}) = (\mathbf{y} - \mathbf{X}\theta)^T(\mathbf{y} - \mathbf{X}\theta) = \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \theta)^2$$

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

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

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

# How to choose the step size ?



# Newton's algorithm

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

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

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

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

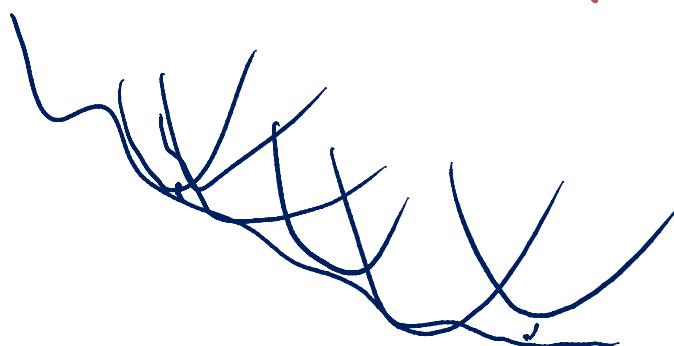
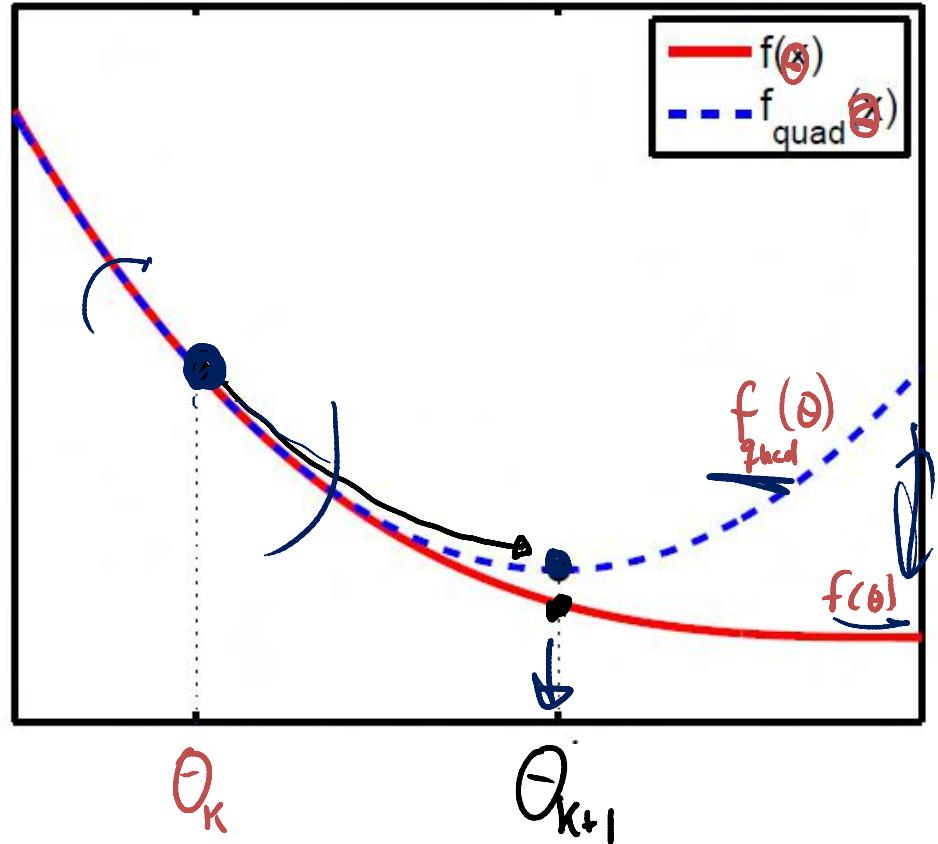
differentiating and equating to zero to solve for  $\theta_{k+1}$ .

$$\nabla f_{quad}(\theta) = \cancel{f(\theta_k)} + \cancel{\mathbf{g}_k} + \mathbf{H}_k(\theta - \theta_k) = 0$$

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

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

# Newton as bound optimization

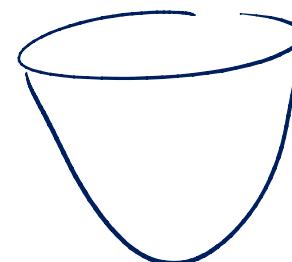


# Newton's algorithm for linear regression

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

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

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



$$\theta_{k+1} = \theta_k - H_k^{-1} \nabla f(\theta_k)$$

$$= \theta_k - (2\mathbf{X}^\top \mathbf{X})^{-1} [-2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\theta_k]$$

$$= \theta_k + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \theta_k$$

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

# 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  $\boxed{\mathbf{H}_k \mathbf{d}_k = -\mathbf{g}_k}$  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      $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \eta_k \mathbf{d}_k$  ✓
-

SGD

epoch

$$\nabla_{\theta} f(\theta) = \underbrace{\int \nabla_{\theta} f(x, \theta) P(x) dx}_{E[\nabla f(x, \theta)]} \stackrel{\text{want}}{=} 0$$

$$\eta \times E[\nabla f(x, \theta)] = 0 \times \eta = 0$$

$$\theta + \eta E[\nabla f(x, \theta)] = \theta$$

$$\theta_{k+1} = \theta_k - \eta E[\nabla f] + \eta [E[\nabla f] - \nabla f_{(x^k)}]$$

$\approx \theta_k - \eta \frac{1}{n} \sum_{i=1}^n \nabla f(x^{(i)}, \theta_k)$

$\approx \theta_k - \eta \nabla f(x^{(k)}, \theta_k)$

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

$\eta \nabla f$

$\theta_k$

$\theta_{k+1}$

$\nabla f_{(x^k)}$

$\approx$

$\text{noise}$

# Online learning with mini-batches

Batch


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

$\left( \begin{array}{l} n \\ \text{data points} \end{array} \right)$

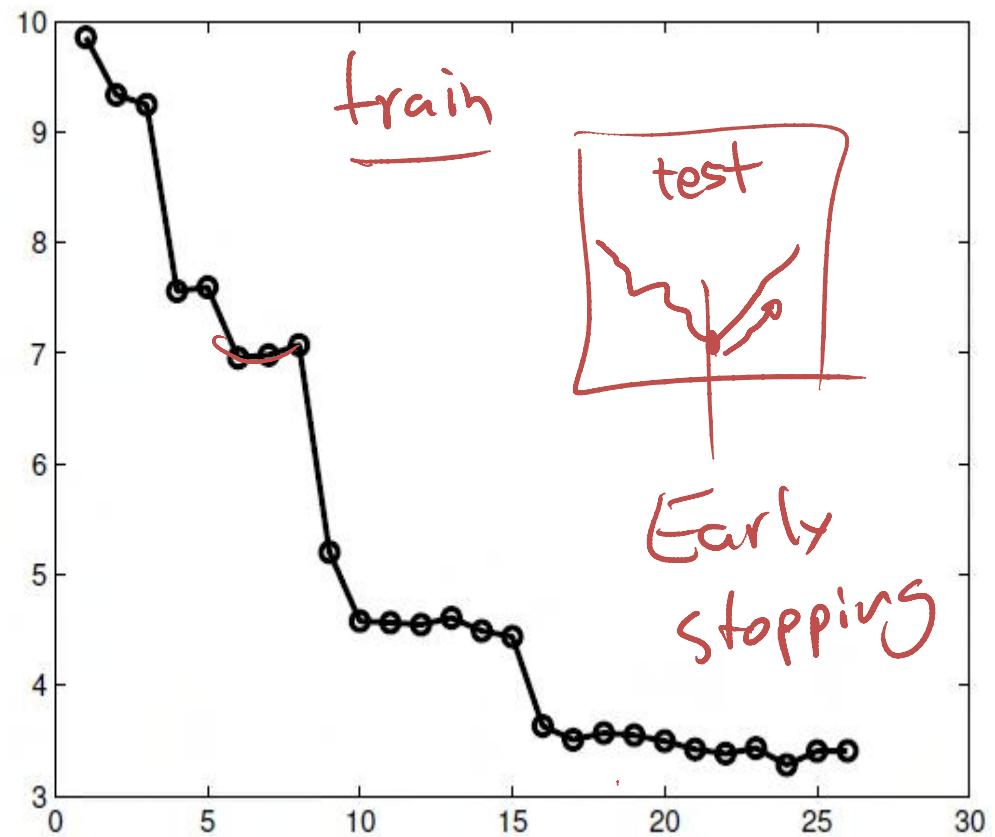
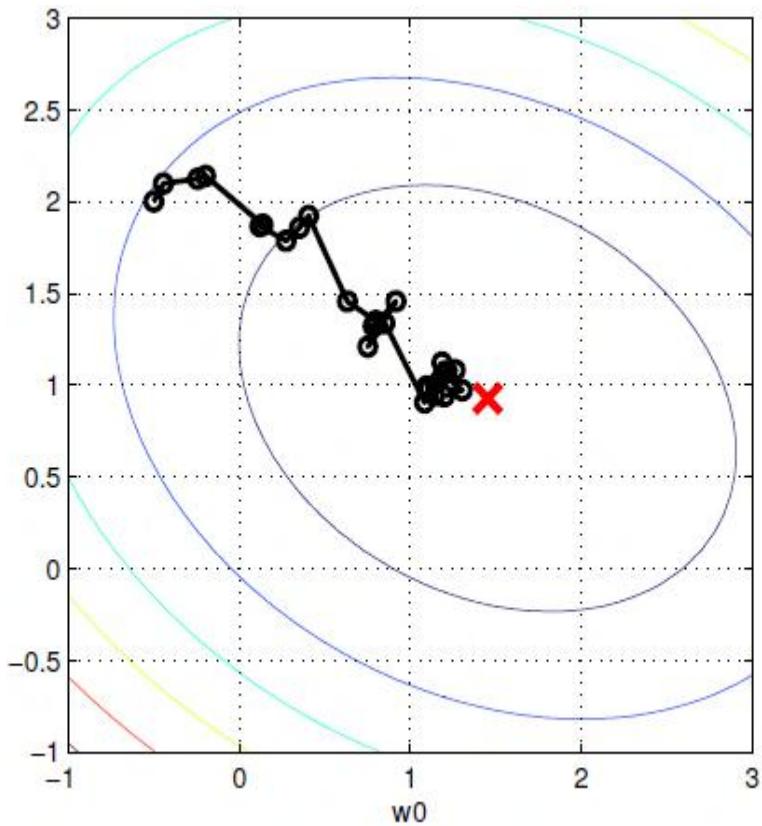
Online

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

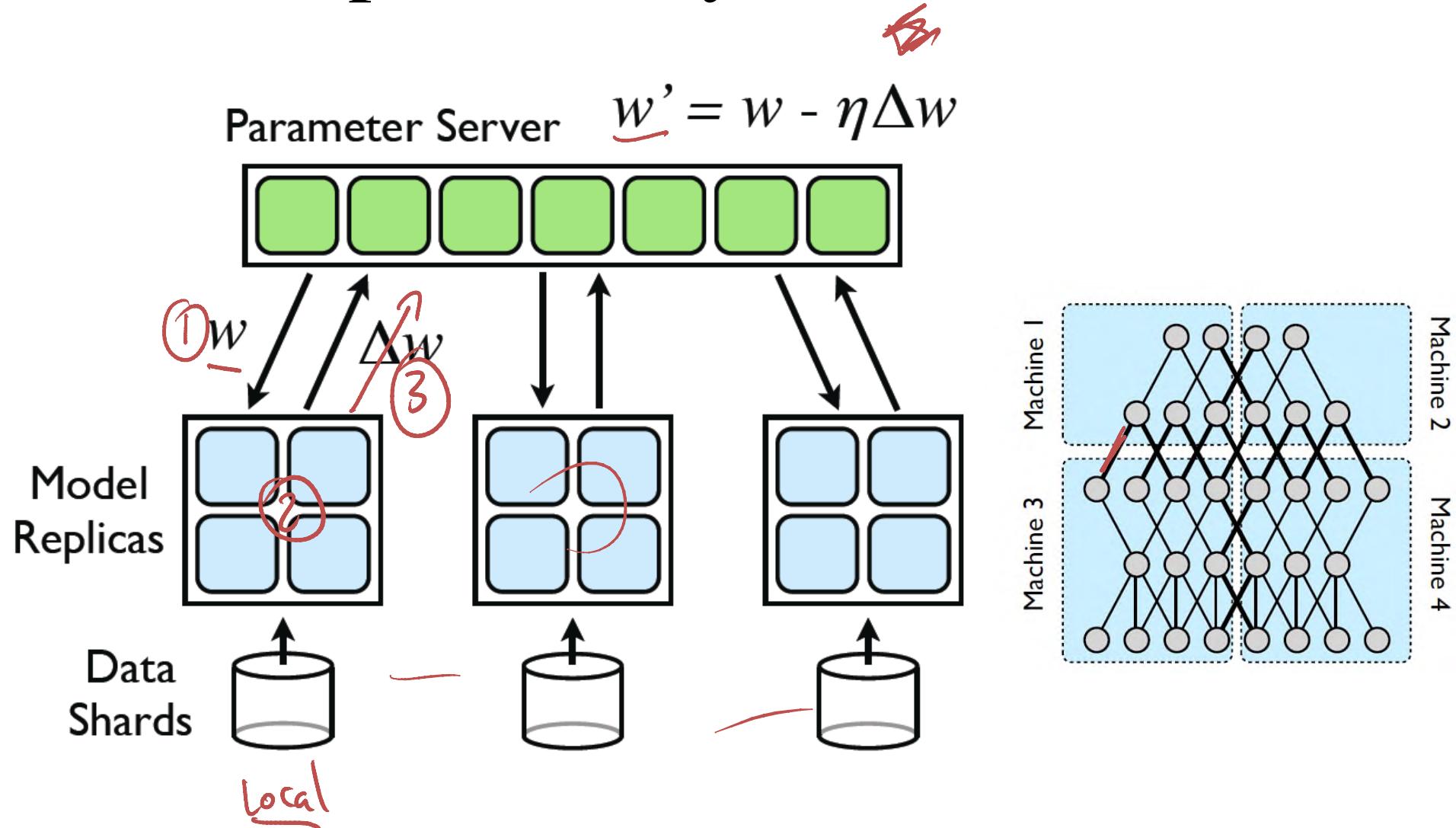
mini-batch


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

# The online learning algorithm



# Downpour – Asynchronous SGD



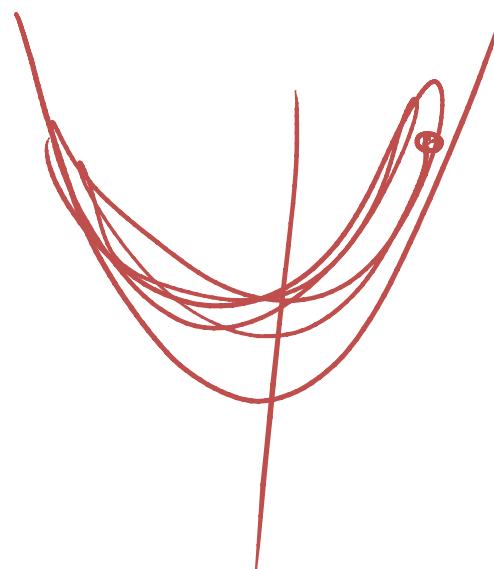
[Jeff Dean et al.]

# Polyak averaging

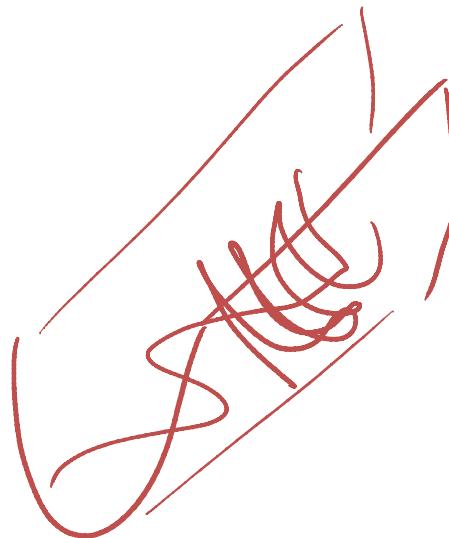
Polyak averaging (see papers of Mark Schmidt / Francis Bach)

$$\frac{1}{t+1} \sum_{i=0}^{t+1} w_i \rightarrow \underbrace{\overline{w}^{(t+1)}}_{\substack{W^{(t+1)} = W^{(t)} - \gamma^{(t)} \nabla L(W^{(t)}, v^{(t)}) \\ \text{---} \\ \overline{W}^{(t)} = \overline{W}^{(t)} - \frac{1}{t} (\overline{W}^{(t)} - W^{(t)})}} = \overline{W}^{(t)} - \frac{1}{t} (\overline{W}^{(t)} - W^{(t)}),$$

See also predictive variance reduction (Tong Zhang, NIPS 2013)



# Momentum



$$\left[ \theta_{t+1} = \theta_t \right] \equiv \alpha(\theta_t - \theta_{t-1}) + (1-\alpha) \left[ -\eta \nabla J(c) \right]$$
$$\Delta \theta_{t+1} \qquad \qquad \Delta \theta_t$$

# Adagrad: Put more weight on rare features

$y_t$	$\phi_{t,1}$	$\phi_{t,2}$	$\phi_{t,3}$
1	1	0	0
-1	.5	0	1
1	-.5	1	0
-1	0	0	0
1	.5	0	0
-1	1	0	0
1	-1	1	0
-1	-.5	0	1

- ① Frequent, irrelevant
- ② Infrequent, predictive
- ③ Infrequent, predictive

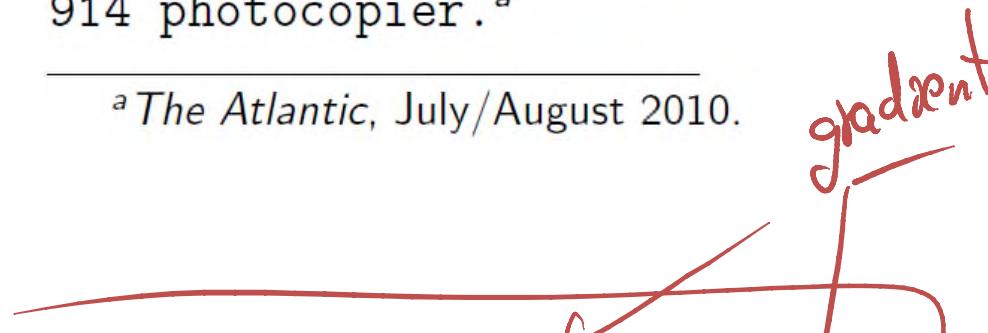
Text data:

The most unsung birthday  
in American business and  
technological history  
this year may be the 50th  
anniversary of the Xerox  
914 photocopier.<sup>a</sup>

<sup>a</sup> *The Atlantic*, July/August 2010.

$$w_i^{(t+1)} \leftarrow w_i^{(t)} - \frac{\eta}{\sqrt{\sum_{\tau=1}^t g_{\tau,i}^2}} g_{t,i}$$

gradient



[Duchi et al.]

# Other useful optimization

- BFGS and limited-BFGS (Take e.g. Nick Trefethen's course)
- Nesterov's method (See Nesterov's book) ✓
- Proximal methods (See Bertsekas book) ✓
- Natural gradient (Yoshua Bengio et al – ICLR) ✓
- Hessian-vector updates and automatic differentiation (Bishop's book and Nocedal & Wright)
- Convex optimization / constrained optimization (See Boy's book)

# Next lecture

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