

Introduction to Machine Learning - CS725  
Instructor: Prof. Ganesh Ramakrishnan  
Lecture 10 - Optimization Foundations Applied to  
Regression Formulations

## Foundations: Necessary Condition 2

- Is  $\nabla^2 f(\mathbf{w}^*)$  positive definite ?  
i.e.  $\forall \mathbf{x} \neq 0$ , is  $\mathbf{x}^T \nabla^2 f(\mathbf{w}^*) \mathbf{x} > 0$ ? (A sufficient condition for local minimum)  
(Any positive definite matrix is also positive semi-definite)  
(Cite : Section 3.12 & 3.12.1)<sup>1</sup>

$$\nabla^2 f(\mathbf{w}^*) = 2\Phi^T \Phi + 2\lambda I \quad (1)$$

$$\implies \mathbf{x}^T \nabla^2 f(\mathbf{w}^*) \mathbf{x} = 2\mathbf{x}^T (\Phi^T \Phi + \lambda I) \mathbf{x} \quad (2)$$

$$= 2 \left( (\Phi + \sqrt{\lambda} I) \mathbf{x} \right)^T \Phi \mathbf{x} \quad (3)$$

$$= 2 \left\| (\Phi + \sqrt{\lambda} I) \mathbf{x} \right\|^2 \geq 0 \quad (4)$$

- And with  $\lambda = 0$ , if  $\Phi$  has full column rank ,

$$\Phi \mathbf{x} = 0 \quad \text{iff} \quad \mathbf{x} = 0 \quad (5)$$

$\therefore$  If  $\mathbf{x} \neq 0$ ,  $\mathbf{x}^T \nabla^2 f(\mathbf{w}^*) \mathbf{x} > 0$

# Example of linearly correlated features

- Example where  $\Phi$  doesn't have a full column rank,

$$\Phi = \begin{bmatrix} x_1 & x_1^2 & x_1^2 & x_1^3 \\ x_2 & x_2^2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ x_n & x_n^2 & x_n^2 & x_n^3 \end{bmatrix} \quad (6)$$

- This is the simplest form of linear correlation of features, and it is not at all desirable.
- Effect of a nonzero  $\lambda$  with such  $\Phi$  is that

# Example of linearly correlated features

- Example where  $\Phi$  doesn't have a full column rank,

$$\Phi = \begin{bmatrix} x_1 & x_1^2 & x_1^2 & x_1^3 \\ x_2 & x_2^2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ x_n & x_n^2 & x_n^2 & x_n^3 \end{bmatrix} \quad (6)$$

- This is the simplest form of linear correlation of features, and it is not at all desirable.
- Effect of a nonzero  $\lambda$  with such  $\Phi$  is that it tends to make the Hessian more positive definite

# Do Closed-form solutions Always Exist?

- Linear regression and Ridge regression both have closed-form solutions

- For linear regression,

$$w^* = (\Phi^T \Phi)^{-1} \Phi^T y$$

- For ridge regression,

$$w^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$$

(for linear regression,  $\lambda = 0$ )

- What about optimizing the formulations (constrained/penalized) of Lasso ( $L_1$  norm)? And support-based penalty ( $L_0$  norm)? **Also requires tools of Optimization/duality**

# Gradient Descent Algorithm

Gradient descent is based on our previous observation that if the multivariate function  $F(\mathbf{x})$  is defined and differentiable in a neighborhood of a point  $\mathbf{a}$ , then  $F(\mathbf{x})$  decreases fastest if one proceeds from  $\mathbf{a}$  in the direction of the negative of the gradient of  $F$  at  $\mathbf{a}$ , i.e.  $-\nabla F(\mathbf{a})$ .

Therefore,

$$\Delta \mathbf{w}^{(k)} = -\nabla \mathbf{E}(\mathbf{w}^{(k)}) \quad (7)$$

Hence,

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + 2\mathbf{t}^{(k)}(\Phi^T \mathbf{y} - \Phi^T \Phi \mathbf{w}^{(k)} - \lambda \mathbf{w}^{(k)}) \quad (8)$$

# Gradient Descent Algorithm

**Find** starting point  $\mathbf{w}^{(0)} \in \mathcal{D}$

- $\Delta \mathbf{w}^k = -\nabla \varepsilon(\mathbf{w}^{(k)})$
- Choose a step size  $t^{(k)} > 0$  using exact or backtracking ray search.
- Obtain  $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + t^{(k)} \Delta \mathbf{w}^{(k)}$ .
- Set  $k = k + 1$ . **until** stopping criterion (such as  $\|\nabla \varepsilon(\mathbf{w}^{(k+1)})\| \leq \epsilon$ ) is satisfied

# Gradient Descent Algorithm

## Exact line search algorithm to find $t^{(k)}$

- The line search approach first finds a descent direction along which the objective function  $f$  will be reduced and then computes a step size that determines how far  $\mathbf{x}$  should move along that direction.
- In general,

$$t^{(k)} = \arg \min_t f(\mathbf{w}^{(k+1)}) \quad (9)$$

- Thus,



# Gradient Descent Algorithm

## Exact line search algorithm to find $t^{(k)}$

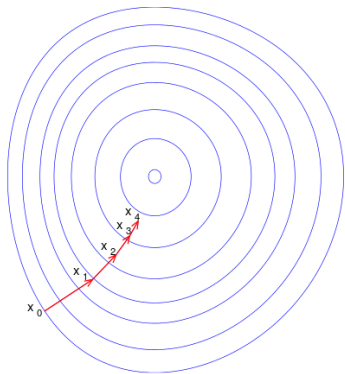
- The line search approach first finds a descent direction along which the objective function  $f$  will be reduced and then computes a step size that determines how far  $\mathbf{x}$  should move along that direction.
- In general,

$$t^{(k)} = \arg \min_t f(\mathbf{w}^{(k+1)}) \quad (9)$$

- Thus,

$$t^{(k)} = \arg \min_t \left( \mathbf{w}^{(k)} + 2t \left( \Phi^T \mathbf{y} - \Phi^T \phi \mathbf{w}^{(k)} - \lambda \mathbf{w}^{(k)} \right) \right) \quad (10)$$

# Example of Gradient Descent Algorithm



**Figure 1:** A red arrow originating at a point shows the direction of the negative gradient at that point. Note that the (negative) gradient at a point is orthogonal to the level curve going through that point. We see that gradient descent leads us to the bottom of the bowl, that is, to the point where the value of the function  $F$  is minimal. Source: Wikipedia

# Constrained Least Squares Linear Regression

Find

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \|\phi \mathbf{w} - \mathbf{y}\|^2 \quad \text{s.t.} \quad \|\mathbf{w}\|_p \leq \zeta, \quad (11)$$

where

$$\|\mathbf{w}\|_p = \left( \sum_{i=1}^n |w_i|^p \right)^{\frac{1}{p}} \quad (12)$$

**Claim: This is an equivalent reformulation of the penalized least squares. Why?**

# p-Norm level curves

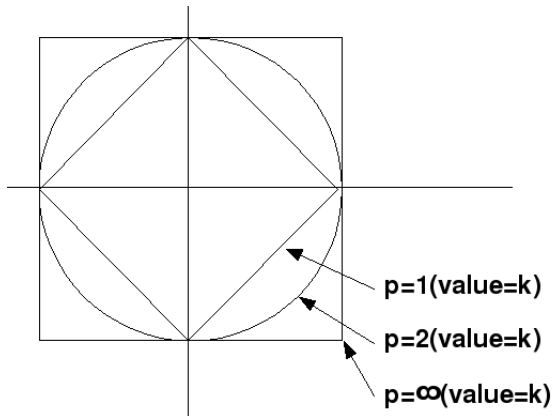


Figure 2: p-Norm curves for constant norm value and different  $p$

# Convex Optimization Problem

- Formally, a convex optimization problem is an optimization problem of the form

$$\text{minimize } f(\mathbf{x}) \quad (13)$$

$$\text{subject to } c \in C \quad (14)$$

where  $f$  is a convex function,  $C$  is a convex set, and  $\mathbf{x}$  is the optimization variable.

- An improved form of the above would be

$$\text{minimize } f(\mathbf{x}) \quad (15)$$

$$\text{subject to } g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m \quad (16)$$

$$h_i(\mathbf{x}) = 0, \quad i = 1, \dots, p \quad (17)$$

where  $f$  is a convex function,  $g_i$  are convex functions, and  $h_i$  are affine functions, and  $\mathbf{x}$  is the vector of optimization variables.

# Constrained convex problems

**Q.** How to solve constrained problems of the above-mentioned type?

**A.** General problem format :

$$\text{Minimize } f(\mathbf{w}) \text{ s.t. } g(\mathbf{w}) \leq 0 \quad (18)$$

