

Introduction to Machine Learning - CS725  
Instructor: Prof. Ganesh Ramakrishnan  
Lecture 4 - Least Squares Linear Regression

# Regression Model

- Training set (this is your data set),  
 $\mathcal{D} = \langle \mathbf{x}_1, \mathbf{y}_1 \rangle, \langle \mathbf{x}_2, \mathbf{y}_2 \rangle, \dots, \langle \mathbf{x}_m, \mathbf{y}_m \rangle$ 
  - Notation (used throughout the course)
  - $m$  = number of training examples
  - $\mathbf{x}$ 's = input variables / features
  - $\mathbf{y}$ 's = output variable "target" variables
  - $(\mathbf{x}, \mathbf{y})$  - single training example
  - $(\mathbf{x}_i, \mathbf{y}_i)$  - specific example ( $i^{\text{th}}$  training example)
  - $i$  is an index to training set
- Need to determine parameters  $\mathbf{w}$  for the function  $f(\mathbf{x}, \mathbf{w})$  which minimizes our error function  $\varepsilon(f(\mathbf{x}, \mathbf{w}), \mathcal{D})$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left\{ \varepsilon(f(\mathbf{x}, \mathbf{w}), \mathcal{D}) \right\}$$

# Linear Regression Model

- Need to determine  $\mathbf{w}$  for the linear function

*specified*  $f(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^p w_i \phi_i(\mathbf{x}_j) = \phi \mathbf{w}$  which minimizes our error function  $\varepsilon(f(\mathbf{x}, \mathbf{w}), \mathcal{D})$  *unspecified*

- $\phi_i$ 's are the basis functions, and let

$$\phi = \begin{bmatrix} \phi_1(\mathbf{x}_1) & \phi_2(\mathbf{x}_1) & \dots & \phi_p(\mathbf{x}_1) \\ \vdots & \vdots & & \vdots \\ \phi_1(\mathbf{x}_m) & \phi_2(\mathbf{x}_m) & \dots & \phi_p(\mathbf{x}_m) \end{bmatrix} \quad (1)$$

- $$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ \vdots \\ y_m \end{bmatrix} \quad (2)$$

# Least Square Linear Regression Model

$$\mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \cdot \\ \cdot \\ w_p \end{bmatrix} \quad (3)$$

$$\mathbf{w}^* = \underset{\mathbf{w}}{\text{arg min}} \left\{ \sum_{j=1}^m \left( \sum_{i=1}^p w_i \phi_i(\mathbf{x}_j) - y_j \right)^2 \right\} \quad (4)$$

Form of errors fn also specified  
Value of error

$$\varepsilon = \min_{\mathbf{w}} \left( \mathbf{w}^T \phi^T \phi \mathbf{w} - 2\mathbf{y}^T \phi \mathbf{w} + \mathbf{y}^T \mathbf{y} \right) \quad (5)$$

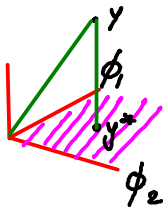
- **Regression**
  - Formal Definition
  - Examples and Types of Regression
- **Least Square Solution**
  - Role of error/loss function
  - Least square solution for linear regression
- **Geometric Interpretation of Least Square Solution**
- **Theorem** :  $\phi^T \phi$  is invertible if and only if  $\phi$  is full column rank

# Geometric Interpretation of Least Square Solution

- Let  $\mathbf{y}^*$  be a solution in the column space of  $\phi$
- The least squares solution is such that the distance between  $\mathbf{y}^*$  and  $\mathbf{y}$  is minimized
- Therefore, the line joining  $\mathbf{y}^*$  to  $\mathbf{y}$  should be orthogonal to the column space

*in column space of  $\phi$*

*to be predicted*



$$\phi \mathbf{w} = \mathbf{y}^* \quad (6)$$

$$(\mathbf{y} - \mathbf{y}^*)^T \phi = 0 \quad (7)$$

$$(\mathbf{y}^*)^T \phi = (\mathbf{y})^T \phi \quad (8)$$

$$(\phi \mathbf{w})^T \phi = \mathbf{y}^T \phi \quad (9)$$

$$\mathbf{w}^T \phi^T \phi = \mathbf{y}^T \phi \quad (10)$$

$$\phi^T \phi \mathbf{w} = \phi^T \mathbf{y} \quad (11)$$

$$\mathbf{w} = (\phi^T \phi)^{-1} \phi^T \mathbf{y} \quad (12)$$

- Here  $\phi^T \phi$  is invertible only if  $\phi$  has full column rank

**Theorem** :  $\phi^T \phi$  is invertible if and only if  $\phi$  is full column rank

Proof :

*only if* Given that  $\phi$  has full column rank and hence columns are linearly independent, we have that  $\phi \mathbf{x} = \mathbf{0} \Rightarrow \mathbf{x} = \mathbf{0}$

Assume on the contrary that  $\phi^T \phi$  is not invertible. Then  $\exists \mathbf{x} \neq \mathbf{0}$  such that  $\phi^T \phi \mathbf{x} = \mathbf{0}$

*Proof by contradiction*

$$\Rightarrow \mathbf{x}^T \phi^T \phi \mathbf{x} = 0$$

$$\Rightarrow (\phi \mathbf{x})^T \phi \mathbf{x} = 0$$

$$\Rightarrow \phi \mathbf{x} = \mathbf{0}$$

*$a^T a = 0 \iff a = 0$*

This is a contradiction. Hence  $\phi^T \phi$  is invertible if  $\phi$  is full column rank

*if* If  $\phi^T \phi$  is invertible then  $\phi \mathbf{x} = \mathbf{0}$  implies  $(\phi^T \phi \mathbf{x}) = \mathbf{0}$ , which in turn implies  $\mathbf{x} = \mathbf{0}$  , **This implies  $\phi$  has full column rank if  $\phi^T \phi$  is invertible. Hence, theorem proved**



- Some more questions on the Least Square Linear Regression Model
- More generally: How to minimize a function?
  - Level Curves and Surfaces
  - Gradient Vector
  - Directional Derivative
  - Hyperplane
  - Tangential Hyperplane
- Gradient Descent Algorithm

# Some questions

If  $A$  is p.d then all  $\lambda(A) > 0 \Rightarrow Ax=0$  has only  $x=0$  as solution, because if  $\exists x \neq 0$  s.t  $Ax=0$  then  $\lambda=0$

will be an eigenvalue of  $A$  ( $\underline{ie}$   $Ax = \lambda x$ )

$\Rightarrow A$  must be invertible ... So suffices to check that  $\phi^T \phi$  has no zero eigenvalues

- What is the relationship between positive definiteness and invertibility?  $\rightarrow x^T(\phi^T \phi)x = \|\phi x\|_2^2 \geq 0 \Rightarrow \forall \lambda(\phi^T \phi) \geq 0$
- When is  $\phi$  not full column rank? What are associated problems and fixes?  $\rightarrow$  eg: if  $m < p$ ,  $\phi_i(x) = \phi_j(x) \forall x$
- How to find a solution if  $\phi$  is not full column rank?

Feature selection

① Select only a subset of  $\phi_i$ 's & drop the rest st subset is linearly independent  
Problem:  $2^n$  subsets to explore!!

$$\sum_{t=1}^k \phi_{i_t}(x) \stackrel{OR}{=} \sum_{t=1}^k \phi_{j_t}(x)$$

Algos for greedily selecting  $\phi_i$  to include or exclude  
eg: Infogain ( $\phi_i$  | attributes selected so far) followed by Least squares soln

② Modify the objective being minimized

Modifying objective?

$$w^k = \underset{w}{\operatorname{argmin}} \sum_{j=1}^n \left[ \left( \sum_{i=1}^p \phi_i(x_j) w_i \right) - y_j \right]^2$$

$$\text{s.t. } \Omega(w) \leq k$$

$$= \underset{w}{\operatorname{argmin}} \|\phi w - y\|^2$$
$$\text{s.t. } \Omega(w) \leq k$$

Eg:  $\Omega(w) = \#$  of non-zero  $w_i$ 's

i.e. Least squares regression s.t. not more than  $k$   $w$ 's are  $\neq 0$  i.e. not more than  $k$   $\phi_i$ 's that are "effective"

→ earlier search space of  $w$



new smaller space of  $w$ 's

3 solns  $\rightarrow$  ① feature selection on  $\phi$ 's to give  $\phi_s$ , followed by  $w^* = (\phi_s^T \phi_s)^{-1} \phi_s^T y$

$$\textcircled{2} w^* = \underset{w}{\operatorname{argmin}} \|\phi w - y\|^2$$

$$\text{s.t. } \Omega(w) \leq \xi$$

$$\textcircled{3} \underbrace{\phi^T \phi} w^* = \phi^T y$$

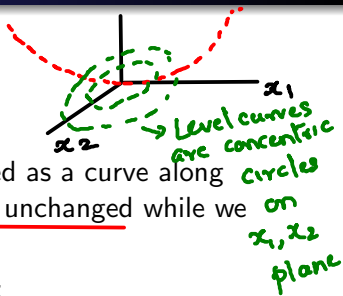
H/w problem: Based on diff inequalities between  $p$  &  $m$ , find cases where this equation has ① No solution ② one solution and ③ multiple solutions

# Solving Least Square Linear Regression Model

- Intuitively: Minimize by setting derivative (gradient) to 0 and find closed form solution.
- For most optimization problems, finding closed form solution is difficult
- Even for linear regression (for which closed form solution exists), are there alternative methods?
- Eg: Consider,  $\mathbf{y} = \phi\mathbf{w}$ , where  $\phi$  is a matrix with full column rank, the least squares solution,  $\mathbf{w}^* = (\phi^T\phi)^{-1}\phi^T\mathbf{y}$ . Now, imagine that  $\phi$  is a very large matrix. with say, 100,000 columns and 1,000,000 rows. Computation of closed form solution might be challenging.
- How about an iterative method?

# Level curves and surfaces

Ex:  $f(x_1, x_2) = x_1^2 + x_2^2$



- A level curve of a function  $f(\mathbf{x})$  is defined as a curve along which the value of the function remains unchanged while we change the value of its argument  $\mathbf{x}$ .
- Formally we can define a level curve as :

$$L_c(f) = \left\{ \underline{\mathbf{x}} \mid f(\mathbf{x}) = c \right\} \quad (13)$$

where  $c$  is a constant.

# Level curves and surfaces

- The image below is an example of different level curves for a single function

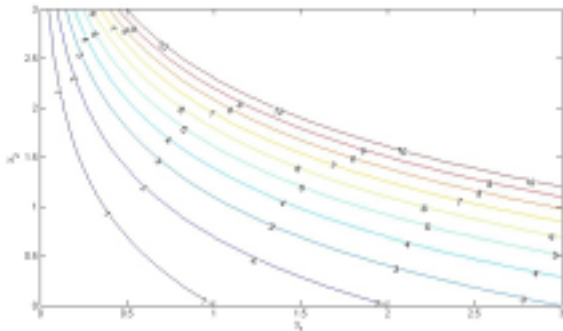


Figure 1: 10 level curves for the function  $f(\mathbf{x}_1, \mathbf{x}_2) = x_1 e^{-x_2}$  (Figure 4.12 from <https://www.cse.iitb.ac.in/~cs709/notes/BasicsOfConvexOptimization.pdf>)

# Directional Derivatives

- Directional derivative: Rate at which the function changes at a given point in a given direction
- The *directional derivative* of a function  $f$  in the direction of a unit vector  $\mathbf{v}$  at a point  $\mathbf{x}$  can be defined as :

$$D_{\mathbf{v}}(f) = \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{v}) - f(\mathbf{x})}{h} \quad (14)$$

Distance Traveled along  $\mathbf{v}$

$$\|\mathbf{v}\| = 1 \quad \mathbf{v} \text{ is treated as a unit vector} \quad (15)$$



# Gradient Vector

$\nabla f$  is direction of maximum directional derivative

- Magnitude (euclidean norm) of gradient vector at any point indicates maximum value of directional derivative at that point
- Direction of gradient vector indicates direction of this maximal directional derivative at that point.
- The *gradient vector* of a function  $f$  at a point  $\mathbf{x}$  is defined as:

$$D_{\mathbf{v}}(f(\mathbf{x})) = \nabla^T f(\mathbf{x}) \mathbf{v}$$



By Cauchy Schwarz Ineq,  $D_{\mathbf{v}}(f(\mathbf{x}))$  is max when  $\mathbf{v}$  is in the direction of  $\nabla f(\mathbf{x})$

$$\nabla f_{\mathbf{x}^*} = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \cdot \\ \cdot \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix} \in \mathbb{R}^n \quad (16)$$

# Gradient Vector

- Magnitude (euclidean norm) of gradient vector at any point indicates maximum value of directional derivative at that point
- The *gradient vector* of a function  $f$  at a point  $\mathbf{x}$  is defined as:

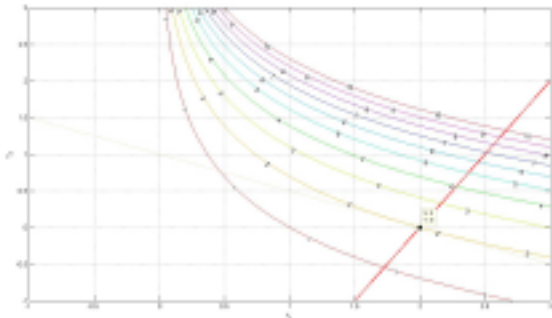
$$\nabla f_{\mathbf{x}^*} = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix} \in \mathbb{R}^n \quad (17)$$

- Thus, at the point of minimum of a differentiable minimization objective (such as least squares for regression), ....

Necessary:  $\nabla \mathcal{E}(w^*) = 0$ . Need to verify that by solving this eqn for least squares regression, we get  $w^* = (\phi^T \phi)^{-1} \phi^T y$

# Gradient Vector

- The figure below gives an example of gradient vector



**Figure 2:** The level curves from Figure 1 along with the gradient vector at  $(2, 0)$ . Note that the gradient vector is perpendicular to the level curve  $x_1 e^{x_2} = 2$  at  $(2, 0)$