

Introduction to Machine Learning - CS725
Instructor: Prof. Ganesh Ramakrishnan
Lecture 9 - Optimization Foundations Applied to
Regression Formulations

Building on questions on Least Squares Linear Regression

- 1 Is there a probabilistic interpretation?
 - Gaussian Error, Maximum Likelihood Estimate
- 2 Addressing overfitting
 - Bayesian and Maximum A posteriori Estimates, Regularization, Support Vector Regression
- 3 How to minimize the resultant and more complex error functions?
 - Level Curves and Surfaces, Gradient Vector, Directional Derivative, Gradient Descent Algorithm, Convexity, Necessary and Sufficient Conditions for Optimality

- 1-norm Error, and L_2 regularized:

- $\min_{w,b,\xi_i,\xi_i^*} \frac{1}{2} \|w\|^2 + C \sum_i (\xi_i + \xi_i^*)$
s.t. $\forall i,$
 $y_i - w^\top \phi(x_i) - b \leq \epsilon + \xi_i,$
 $b + w^\top \phi(x_i) - y_i \leq \epsilon + \xi_i^*,$
 $\xi_i, \xi_i^* \geq 0$

- 2-norm Error, and L_2 regularized:

- $\min_{w,b,\xi_i,\xi_i^*} \frac{1}{2} \|w\|^2 + C \sum_i (\xi_i^2 + \xi_i^{*2})$
s.t. $\forall i,$
 $y_i - w^\top \phi(x_i) - b \leq \epsilon + \xi_i,$
 $b + w^\top \phi(x_i) - y_i \leq \epsilon + \xi_i^*$
- Here, the constraints $\xi_i, \xi_i^* \geq 0$ are not necessary

Need for Optimization so far

- **Unconstrained (Penalized) Optimization:**

$$\mathbf{w}_{Reg} = \arg \min_{\mathbf{w}} \|\phi \mathbf{w} - \mathbf{y}\|_2^2 + \Omega(\mathbf{w})$$

- **Constrained Optimization 1:**

$$\mathbf{w}_{Reg} = \arg \min_{\mathbf{w}} \|\phi \mathbf{w} - \mathbf{y}\|_2^2$$

such that $\Omega(\mathbf{w}) \leq \theta$

- **Constrained Optimization 2 ($t = 1$ or 2):**

$$\arg \min_{w, b, \xi_i, \xi_i^*} \frac{1}{2} \|w\|^2 + C \sum_i (\xi_i^t + \xi_i^{*t})$$

s.t. $\forall i, y_i - w^\top \phi(x_i) - b \leq \epsilon + \xi_i$; $b + w^\top \phi(x_i) - y_i \leq \epsilon + \xi_i^*$

- **Equivalence:** λ (Penalized) $\equiv \theta$ (Constrained)
- **Duality:** Dual of Support Vector Regression

Solving Unconstrained Minimization Problem

- Intuitively: Minimize by setting derivative (gradient) to 0 and hoping to find **closed form** solution.
- When is such a solution a global minimum?
- For most optimization problems, finding closed form solutions is difficult. Even for linear regression (for which closed form solution exists), are there alternative methods?
 - Eg: Consider, $\mathbf{y} = \phi\mathbf{w}$, where ϕ is a matrix with full column rank, the least squares solution, $\mathbf{w}^* = (\phi^T\phi)^{-1}\phi^T\mathbf{y}$. Now, imagine that ϕ 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 iterative methods?

Foundations: Level curves and surfaces

- A level curve of a function $\mathbf{f}(\mathbf{x})$ 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 :

$$L_c(\mathbf{f}) = \left\{ \mathbf{x} \mid \mathbf{f}(\mathbf{x}) = \mathbf{c} \right\} \quad (1)$$

where c is a constant.

Foundations: Level curves and surfaces

- Example of different level curves for a single function

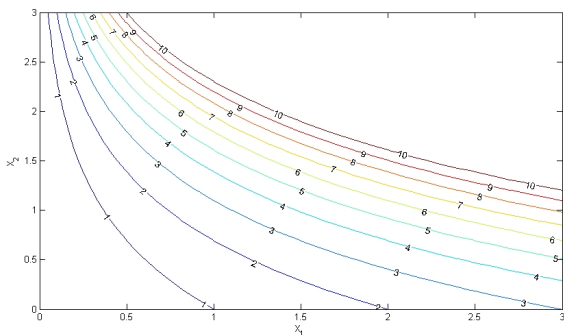


Figure 1: 10 level curves for the function $f(\mathbf{x}_1, \mathbf{x}_2) = x_1 e^{x_2}$ (Figure 4.12 from <https://www.cse.iitb.ac.in/~CS725/notes/classNotes/BasicsOfConvexOptimization.pdf>)

Foundations: Directional Derivatives

- Directional derivative: Rate at which the function changes at a given point \mathbf{x} in a given direction \mathbf{v}
- 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 :

$$D_{\mathbf{v}}(f, \mathbf{x}) = \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{v}) - f(\mathbf{x})}{h} \quad (2)$$

$$\text{s.t. } \|\mathbf{v}\|_2 = 1 \quad (3)$$

- The **gradient vector** of a function f at a point \mathbf{x} is defined as:

$$\nabla f_{\mathbf{x}^*} = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \cdot \\ \cdot \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix} \in \mathbb{R}^n \quad (4)$$

- **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.

Foundations: Gradient Vector

- The figure below illustrates the gradient vector for the same level curves

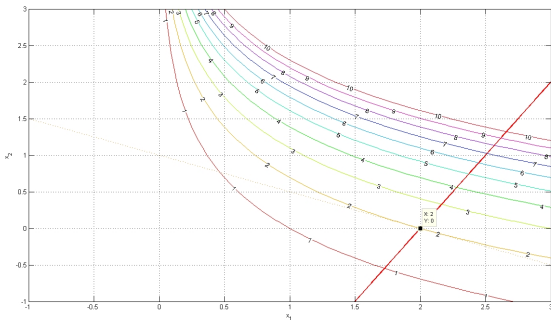


Figure 2: The level curves along with the gradient vector at $(2, 0)$. Note that the gradient vector is perpendicular to the level curve $x_1 e^{x_2} = 2$ at $(2, 0)$

Hyperplanes

- A hyperplane in an n -dimensional Euclidean space is a flat, $n-1$ dimensional subset of that space that divides the space into two disjoint half-spaces.
- Technically, a hyperplane is a set of points whose direction *w.r.t.* a point \mathbf{q} is orthogonal to a vector \mathbf{v} :

$$H_{\mathbf{v},\mathbf{q}} = \left\{ \mathbf{p} \mid (\mathbf{p} - \mathbf{q})^T \mathbf{v} = 0 \right\} \quad (5)$$

- **Tangential Hyperplane:** Plane orthogonal to the gradient vector at \mathbf{x}^* .

Hyperplanes

- A hyperplane in an n -dimensional Euclidean space is a flat, $n-1$ dimensional subset of that space that divides the space into two disjoint half-spaces.
- Technically, a hyperplane is a set of points whose direction *w.r.t.* a point \mathbf{q} is orthogonal to a vector \mathbf{v} :

$$H_{\mathbf{v},\mathbf{q}} = \left\{ \mathbf{p} \mid (\mathbf{p} - \mathbf{q})^T \mathbf{v} = 0 \right\} \quad (5)$$

- **Tangential Hyperplane:** Plane orthogonal to the gradient vector at \mathbf{x}^* .

$$TH_{\mathbf{x}^*} = \left\{ \mathbf{p} \mid (\mathbf{p} - \mathbf{x}^*)^T \nabla f(\mathbf{x}^*) = 0 \right\} \quad (6)$$

We recall that the problem was to find \mathbf{w} such that

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \|\phi \mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|^2 \quad (7)$$

$$= \arg \min_{\mathbf{w}} (\mathbf{w}^T \phi^T \phi \mathbf{w} - 2\mathbf{w}^T \phi \mathbf{y} - \mathbf{y}^T \mathbf{y} + \lambda \|\mathbf{w}\|^2) \quad (8)$$

Foundations: Gradient Vector

- Magnitude (euclidean norm) of gradient vector at any point indicates maximum value of directional derivative at that point
- Thus, at the point of minimum of a differentiable minimization objective (such as least squares for regression),

Foundations: Necessary condition 1

- If $\nabla f(\mathbf{w}^*)$ is defined & \mathbf{w}^* is local minimum/maximum, then $\nabla f(\mathbf{w}^*) = 0$ (A necessary condition) (Cite : Theorem 60) of `CS725/notes/classNotes/BasicsOfConvexOptimization.pdf`
- Given that

$$f(\mathbf{w}) = \arg \min_{\mathbf{w}} (\mathbf{w}^T \phi^T \phi \mathbf{w} - 2\mathbf{w}^T \phi^T \mathbf{y} - \mathbf{y}^T \mathbf{y} + \lambda \|\mathbf{w}\|^2)$$

\implies

- We would have

.....

\implies

\implies

Foundations: Necessary condition 1

- If $\nabla f(\mathbf{w}^*)$ is defined & \mathbf{w}^* is local minimum/maximum, then $\nabla f(\mathbf{w}^*) = 0$ (A necessary condition) (Cite : Theorem 60) [CS725/notes/classNotes/BasicsOfConvexOptimization.pdf](#)
- Given that

$$f(\mathbf{w}) = \arg \min_{\mathbf{w}} (\mathbf{w}^T \phi^T \phi \mathbf{w} - 2\mathbf{w}^T \phi^T \mathbf{y} - \mathbf{y}^T \mathbf{y} + \lambda \|\mathbf{w}\|_2^2)$$
$$\implies \nabla f(\mathbf{w}) = 2\phi^T \phi \mathbf{w} - 2\phi^T \mathbf{y} + 2\lambda \mathbf{w} \quad (10)$$

- We would have

$$\nabla f(\mathbf{w}^*) = 0 \quad (11)$$

$$\implies 2(\phi^T \phi + \lambda I)\mathbf{w}^* - 2\phi^T \mathbf{y} = 0 \quad (12)$$

$$\implies \mathbf{w}^* = (\phi^T \phi + \lambda I)^{-1} \phi^T \mathbf{y} \quad (13)$$

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)
(Note : Any positive definite matrix is also positive semi-definite)
(Cite : Section 3.12 & 3.12.1)¹

.....
 \implies
.....
.....

- And if ϕ has full column rank ,

.....

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

¹CS725/notes/classNotes/LinearAlgebra.pdf

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)²

$$\nabla^2 f(\mathbf{w}^*) = 2\phi^T \phi + 2\lambda I \quad (14)$$

$$\implies \mathbf{x}^T \nabla^2 f(\mathbf{w}^*) \mathbf{x} = 2\mathbf{x}^T (\phi^T \phi + \lambda I) \mathbf{x} \quad (15)$$

$$= 2 \left((\phi + \sqrt{\lambda} I) \mathbf{x} \right)^T \phi \mathbf{x} \quad (16)$$

$$= 2 \left\| (\phi + \sqrt{\lambda} I) \mathbf{x} \right\|^2 \geq 0 \quad (17)$$

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

$$\phi \mathbf{x} = 0 \quad \text{iff} \quad \mathbf{x} = 0 \quad (18)$$

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

Example of linearly correlated features

- Example where ϕ doesn't have a full column rank,

$$\phi = \begin{bmatrix} x_1 & x_1^2 & x_1^2 & x_1^3 \\ x_2 & x_2^2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ x_n & x_n^2 & x_n^2 & x_n^3 \end{bmatrix} \quad (19)$$

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

Example of linearly correlated features

- Example where ϕ doesn't have a full column rank,

$$\phi = \begin{bmatrix} x_1 & x_1^2 & x_1^2 & x_1^3 \\ x_2 & x_2^2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ x_n & x_n^2 & x_n^2 & x_n^3 \end{bmatrix} \quad (19)$$

- This is the simplest form of linear correlation of features, and it is not at all desirable.
- Effect of a nonzero λ with such ϕ 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**