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

1. Is there a probabilistic interpretation?
   - Gaussian Error, Maximum Likelihood Estimate
2. Addressing overfitting
   - Bayesian and Maximum Aposteriori 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

KKT Conditions
① Dual of SVR.. Kernels
② Equivalence of penalized & constrained regression

- 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. $\forall i$,
    $y_i - w^\top \phi(x_i) - b \leq \epsilon + \xi_i,$ }
    $b + w^\top \phi(x_i) - y_i \leq \epsilon + \xi_i^*,$ } Number of constraints
    $\xi_i, \xi_i^* \geq 0$  $= 2 \cdot$ # of examples (m)

- 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

- **Unconstrained (<span style="color:red">Penalized</span>) Optimization:**

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

- **<span style="color:brown">Constrained</span> Optimization 1:**

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

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

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

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

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

- **Equivalence:** $\lambda$ (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?

  $w^*$ s.t $\nabla_f = 0$

  - 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.
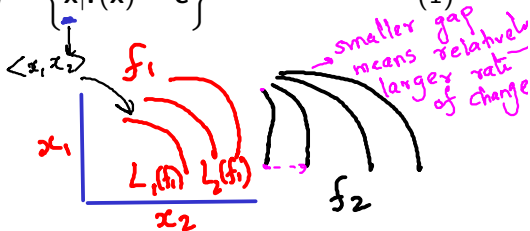
- How about iterative methods?

  $w^{new} = w^{old} + \Delta w$

- 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 x.

- Formally we can define a level curve as :

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

where c is a constant.

- Example of different level curves for a single function
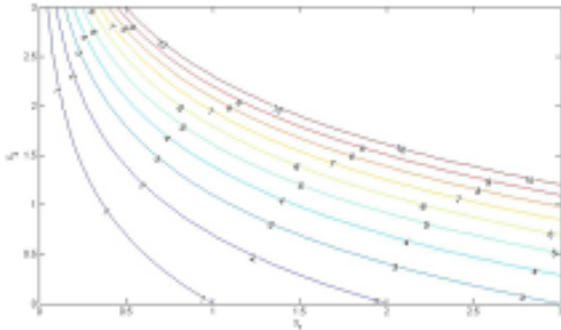


Figure 1: 10 level curves for the function $f(x_1, x_2) = x_1 e^{x_2}$ (Figure 4.12 from https://www.cse.iitb.ac.in/~CS725/notes/classNotes/BasicsOfConvexOptimization.pdf)

- Directional derivative: Rate at which the function changes at a given point **x** in a given direction **v**
- The *directional derivative* of a function $f$ in the direction of a unit vector **v** at a point **x** can be defined as :
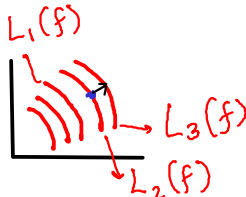
$$D_{\mathbf{v}}(f, \mathbf{x}) = \lim_{h \to 0} \frac{f(\mathbf{x} + h\mathbf{v}) - \mathbf{f}(\mathbf{x})}{h} \qquad (2)$$

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

Claim: $D_v(f, x) = v^T \nabla f(x)$

- The **g**radient vector of a function $f$ at a point **x** is defined as:

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

$L_1(f)$

$L_3(f)$

$L_2(f)$

$\|\nabla f_x\|_2$

- 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.

$\dfrac{\nabla f_x}{\|\nabla f_x\|}$

# Foundations: Gradient Vector

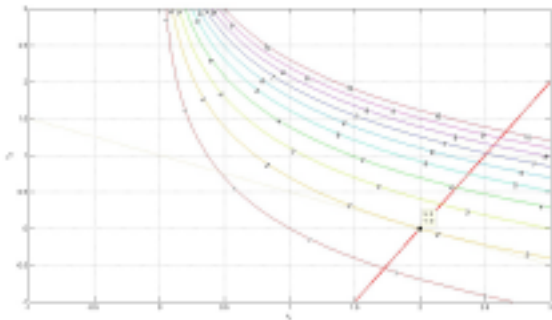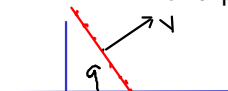- The figure below illustrates the gradient vector for the same level curves



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 $\mathbf{v}$:

$$H_{\mathbf{v},\mathbf{q}} = \left\{ \mathbf{p} \,\middle|\, (\mathbf{p} - \mathbf{q})^{\mathsf{T}}\mathbf{v} = \mathbf{0} \right\} \qquad (5)$$

$$\{p \mid (p-q)^{\mathsf{T}} v = 0\} \qquad \text{such that} \qquad p, q, v \in \mathbb{R}^n$$

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

$$TH_{x^*} = H_{\nabla f(x^*), x^*} = \{p \mid (p - x^*)^{\mathsf{T}} \nabla f(x^*) = 0\}$$

# 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 $\mathbf{v}$:

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

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

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

We recall that the problem was to find **w** such that $\left(L_2 \text{ regularized linear regression}\right)$

$$
\begin{aligned}
\mathbf{w}^* &= \arg\min_{\mathbf{w}} \|\phi\mathbf{w} - \mathbf{y}\|^2 + \lambda\|\mathbf{w}\|^2 \qquad (7) \\
&= \arg\min_{\mathbf{w}}(\mathbf{w}^T\phi^T\phi\mathbf{w} - 2\mathbf{w}^T\phi\mathbf{y} + \mathbf{y}^T\mathbf{y} + \lambda\|\mathbf{w}\|^2) \quad (8)
\end{aligned}
$$

- 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), ....

$$\text{We expect} \quad \nabla f(\omega^*) = 0$$

- 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

$$\text{Quadratic in } \omega \cdots \frac{d\phi^2\omega^2}{d\omega} = 2\phi^2\omega$$

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

$$\implies \ldots\ldots\ldots \nabla f(\omega) = \begin{bmatrix} \frac{\partial f_{\omega_1}}{} \\ \frac{\partial f_{\omega_2}}{} \\ \vdots \\ \frac{\partial f_{\omega_n}}{} \end{bmatrix} = \begin{matrix} \nabla_\omega(\omega^T\phi^T\phi\omega) - \nabla_\omega(2\omega^T\phi y) \\ + \nabla_\omega(\lambda\omega^T\omega) \end{matrix}$$

- We would have

$$= 2\phi^T\phi\omega - 2\underline{\phi y} + 2\lambda\omega$$

$$\ldots\ldots\ldots$$

$$\implies \ldots\ldots\ldots\ldots\ldots$$

$$\implies \ldots\ldots\ldots\ldots\ldots$$

- 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

$$f(\mathbf{w}) = \underset{\mathbf{w}}{\arg\min}(\mathbf{w}^T\phi^T\phi\mathbf{w} - 2\mathbf{w}^T\phi^T\mathbf{y} - \mathbf{y}^T\mathbf{y} + \lambda\|\mathbf{w}\|^2) \quad (9)$$

$$\implies \nabla f(\mathbf{w}) = 2\underbrace{\phi^T\phi\mathbf{w} - 2\phi^T\mathbf{y}}_{\text{Disappears at } \mathbf{w}^*} + 2\lambda\mathbf{w} \quad (10)$$

- We would have

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

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

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

$2\lambda w = 2\lambda I w \cdots \text{Different representation}$

$\text{Assuming invertibility}$

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

*(handwritten, top right)* Necessary for Local min that $\nabla^2 f(x^*) \succeq 0$

$$\nabla^2 f(\omega) = \begin{bmatrix} \dfrac{\partial^2 f(\omega)}{\partial \omega_i^2} & \cdots \\ & \dfrac{\partial^2 f(\omega)}{\partial \omega_i \omega_j} \end{bmatrix}$$

Hessian is symmetric

$$\nabla f(\omega) = 2(\phi^T \phi + \lambda I)\omega - 2\phi^T y$$

$$\nabla^2 f(\omega) = 2(\phi^T \phi + \lambda I)$$

*(handwritten)* can be ignored for $\nabla^2 f(\omega)$

$\forall v \neq 0 \dots v^T \nabla^2 f(\omega) v \geq 0$

Because $\dots (\phi v + \sqrt{\lambda} v)^T (\phi v + \sqrt{\lambda} v)$

- And if $\phi$ **has full column rank** , $= \|\phi v + \sqrt{\lambda} v\|_2^2 \geq 0$

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

# Foundations: Necessary Condition 2

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

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

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

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

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

(About positive semidefiniteness)

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

$$\phi \mathbf{x} = 0 \quad iff \quad \mathbf{x} = 0 \tag{18}$$

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

[2]CS725/notes/classNotes/LinearAlgebra.pdf

$\rightarrow \|p\| = 0$ iff $p \leq 0$
i.e. $(\phi + \sqrt{\lambda}I)x = 0$
$\Leftrightarrow x^T \nabla^2 f(\omega^*)x = 0$
$x^T \nabla^2 f(\omega^*)x = 0$ iff $x = 0$
$\Rightarrow \nabla^2 f(\omega^*) \geq 0$

① with any $\lambda \geq 0$ $\nabla^2 f(\omega) \geq 0$ $\forall \omega \in \mathbb{R}^n$

The ridge regression
objective function is convex
==everywhere== (& ∴ also convex at $\omega^*$)

global! min



$f(\omega)$

$\nabla^2 f(\omega) \geq 0$

$\omega^*$ $\omega$

Necessary condition for
local min at $\omega^*$ is
"convex" or cup-shaped
curvature at $\omega^*$

& $\Phi$ full col·rank

② with $\lambda = 0$ $\wedge$ $\nabla^2 f(\omega) > 0$

The ridge regression
objective function is strictly convex
==everywhere== (& ∴ also at $\omega^*$)

global min



$f(\omega)$

$\omega$

Sufficient condition for
local (& even global) min at $\omega^*$
is strictly convex curvature

New takeaways:

① $\nabla^2 f(\omega) \geq 0 \quad \forall \omega \Rightarrow f$ is convex
everywhere &

∴ necessary condition
for local min to become
global min

② $\nabla^2 f(\omega) > 0 \quad \forall \omega \Rightarrow f$ is strictly convex
everywhere &

sufficient condition for
local min to become
global min

③ If $\lambda > 0$, $\nabla^2 f(\omega)$ tends to become "more"
positive definite

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

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

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

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

- 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,

$$w^* = (\phi^\top \phi)^{-1} \phi^\top y$$

  - For ridge regression,

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

  (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