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*

## SVR objective

- 1-norm Error, and $L_2$ regularized:
  - $\min_{\mathbf{w}, b, \xi_i, \xi_i^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*)$
    s.t. $\forall i,$
    $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon + \xi_i,$
    $b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i \leq \epsilon + \xi_i^*,$ $\left.\right\}$ Number of constraints = 2· # of examples (m)
    $\xi_i, \xi_i^* \geq 0$

- 2-norm Error, and $L_2$ regularized:
  - $\min_{\mathbf{w}, b, \xi_i, \xi_i^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i^2 + \xi_i^{*2})$
    s.t. $\forall i,$
    $y_i - \mathbf{w}^\top \phi(x_i) - b \leq \epsilon + \xi_i,$
    $b + \mathbf{w}^\top \phi(x_i) - y_i \leq \epsilon + \xi_i^*$
  - Here, the constraints $\xi_i, \xi_i^* \geq 0$ are not necessary

- **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 $2$):**

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

s.t. $\forall i,\ y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon + \xi_i;\ b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i \leq \epsilon + \xi_i^*$

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

- 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?

  $$w^* \text{ s.t } \nabla_f = 0$$

  - 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.
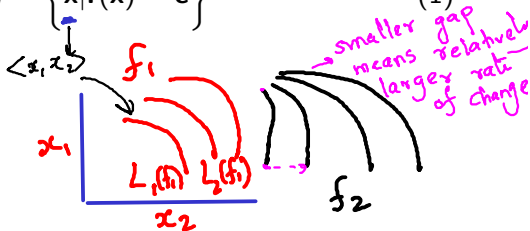
- How about iterative methods?

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

- 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.

# Foundations: Level curves and surfaces
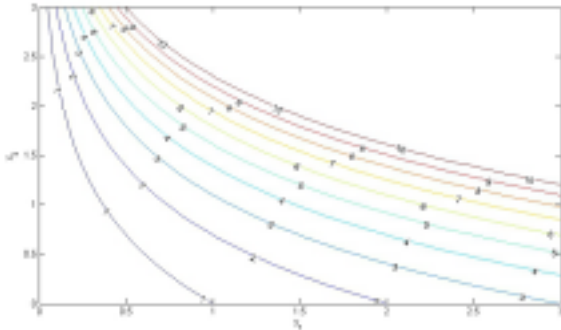
- Example of different level curves for a single function



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 :

$$D_{\mathbf{v}}(f, \mathbf{x}) = \lim_{h \to 0} \frac{f(\mathbf{x} + h\mathbf{v}) - \mathbf{f}(\mathbf{x})}{h} \qquad (2)$$
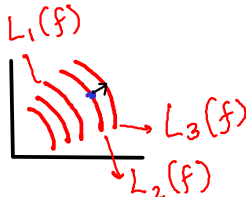
$$s.t. \ ||\mathbf{v}||_2 = \mathbf{1} \qquad (3)$$

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:



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

$\|\nabla f_x\|_2$

- 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.

$\left.\vphantom{\begin{matrix}a\\b\end{matrix}}\right\} \dfrac{\nabla f_x}{\|\nabla f_x\|}$

# Foundations: Gradient Vector

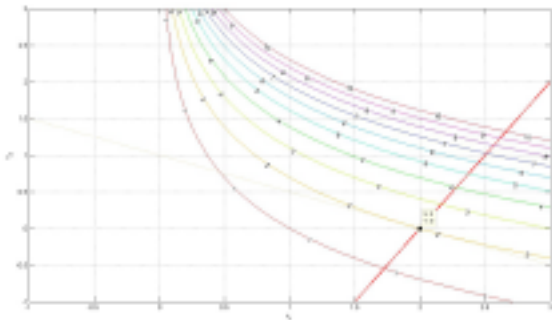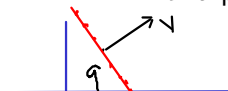- The figure below illustrates the gradient vector for the same level curves



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)

- 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 **q** is orthogonal to a vector **v**:

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

such that $\quad p, q, v \in \mathbb{R}^n$

$\{ p \mid (p-q)^{\mathsf{T}} v = 0 \}$

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

$$TH_{x^*} = H_{\nabla f(x^*), x^*} = \left\{ p \;\middle|\; (p - x^*)^{\mathsf{T}} \nabla f(x^*) = 0 \right\}$$

$\nabla f(x^*)$

$x^*$

## 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 $\mathbf{v}$:

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

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

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

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

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

- 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), ....

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

- 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

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

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

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

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

- We would have

$$\cdots\cdots\cdots$$

$$\implies \cdots\cdots\cdots\cdots\cdots$$

$$\implies \cdots\cdots\cdots\cdots\cdots$$

- 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

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

$$\implies \nabla f(\mathbf{w}) = 2\Phi^T \Phi \mathbf{w} - 2\Phi^T \mathbf{y} + 2\lambda \mathbf{w} \quad (10)$$

*Disappears at $\mathbf{w}^*$*

- We would have

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

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

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

*$2\lambda w = 2\lambda I w$ .. Different representation*

*Assuming invertibility*

→ Necessary for Local min that $\nabla^2 f(x^*) \succeq 0$

- 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]

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

Hessian is symmetric

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

$\implies \nabla^2 f(\omega) = 2(\phi^T \phi + \lambda I)$    can be ignored for $\nabla^2 f(\omega)$

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

Because $2\|\phi v\|_2^2 + \lambda \|v\|_2^2 \geq 0$

- And if Φ **has full column rank**,

...............

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

- 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** , $\|P\| = 0$ iff $P \leq 0$

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

$\stackrel{ie}{=} (\phi + \sqrt{\lambda}I)x = 0$

$\Leftrightarrow x^T \nabla^2 f(\omega^*)x = 0$

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

$x^T \nabla^2 f(\omega^*)x = 0$ iff $x \leq 0$

$\Rightarrow \nabla^2 f(\omega^*) \succ 0$

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

$\Downarrow$

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

_global! min_



$f(\omega)$

$\omega^*$   $\nabla^2 f(\omega) \succeq 0$   $\rightarrow \omega$

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

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

& $\phi$ full col. rank

$\Downarrow$

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

_global min_



$f(\omega)$

$\rightarrow \omega$

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 becom 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

- Linear regression and Ridge regression both have closed-form solutions
  - For linear regression,

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

  - For ridge regression,

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

  (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