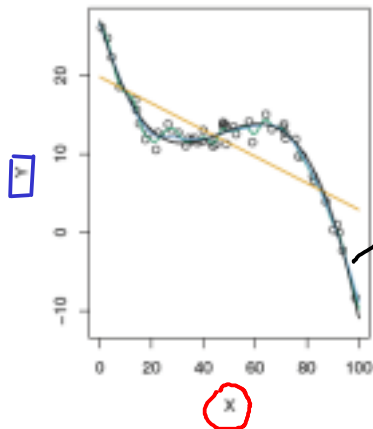Introduction to Machine Learning - CS725
Instructor: Prof. Ganesh Ramakrishnan
Lecture 8 - Support Vector Regression and
Optimization Basics

# Recap: Overfitting and Regularization through Illustration



$$y \simeq w_0 + w_1 x + w_2 x^2 + w_3 x^3$$

- Consider a degree 3 polynomial regression model as shown in the figure
- Each bend in the curve corresponds to increase in $\|w\|$
- Eigen values of $(\Phi^\top \Phi + \lambda I)$ are indicative of curvature. Increasing $\lambda$ reduces the curvature [Prob 8 of Tut 3 & 4]

# Support Vector Regression

One more formulation before we look at Tools of Optimization/duality

# Building on questions on Least Squares Linear Regression

1. Is there a probabilistic interpretation?
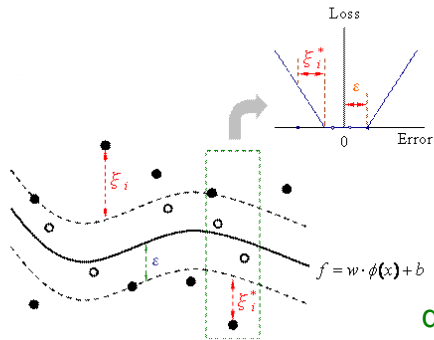   - Gaussian Error, Maximum Likelihood Estimate
2. Addressing overfitting
   - Bayesian and Maximum Aposteriori Estimates, Regularization, **Support Vector Regression**

   $L(w, x, y) \leftarrow$

   $\frac{1}{2}(\|w\|)$

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

# Support Vector Regression (SVR)



SVR attempts to avoid overfitting also through the loss component by introducing an epsilon insensitive loss (band)

$$\min_{\omega} \; L_{\epsilon}(\omega, D) + \Omega(\omega)$$

discover the regression curve s.t points around it within epsilon band have not been penalized

- $\epsilon$-*insensitive loss*: Any point in the band (of $\epsilon$) is not penalized.
- Any point outside the band is penalized, and has slackness $\xi_i$ or $\xi_i^*$  SVR is graceful on epsilon measurement errors and penalizes rest
- The SVR model curve may not pass through any training point

$$\text{Loss}(x_i) \geq 0$$

$$\left(\xi_i\right) = y_i - \left(\omega^T \phi(x_i) + b + \epsilon\right) \quad \text{if} \quad y_i \geq \omega^T \phi(x_i) + b + \epsilon$$

$$\left(\xi_i^{\rightarrow}\right) = \omega^T \phi(x_i) + b - \epsilon - y_i \quad \text{if} \quad y_i \leq \omega^T \phi(x_i) + b - \epsilon$$

Need a compact(er) expression for the Loss!!

$$\text{Loss} = \sum_i \text{Loss}(x_i) = \sum_i \left(\xi_i + \xi_i^*\right) \quad \begin{bmatrix} \text{since at most one} \\ \text{of them should} \\ \text{be non-zero} \end{bmatrix}$$

<u>INTENT</u>

$$\xi_i = \max\left(0, \, y_i - \left(\omega^T \phi(x_i) + b + \epsilon\right)\right)$$

$$\xi_i^{\rightarrow} = \max\left(0, \, \omega^T \phi(x_i) + b - \epsilon - y_i\right)$$

Execution: $\min a \quad s.t \quad a = \max(b, c) \Leftrightarrow \min a \quad s.t \quad \begin{array}{c} a \geq b \\ a \geq c \end{array}$

- The tolerance $\epsilon$ is fixed
- It is desirable that $\forall i$:

(a) $\xi_i = y_i - (w^T\phi(x_i) + b + \epsilon)$ → Regression curve (line)

(b) $\xi_i = 0$
} $+\epsilon$ band

(c) $\xi_i^* = 0$
} $-\epsilon$ band

(d) $\xi_i^* = w^T\phi(x_i) + b - \epsilon - y_i$

- The tolerance $\epsilon$ is fixed
- It is desirable that $\forall i$: [Necessary conditions that become sufficient when also
  - $y_i - w^\top\phi(x_i) - b \leq \epsilon + \xi_i$
  - $b + w^\top\phi(x_i) - y_i \leq \epsilon + \xi_i^*$

  minimize $\sum_i \left(\xi_i + \xi_i^*\right)$]

$0 \leq \xi_i, \xi_i^*$

obvious from constr -uction → $\xi_i \cdot \xi_i^* = 0$ → Claim: This is redundant constraint

$\epsilon$ is a hyperparameter like $\lambda$ or $\sigma^2$

**Claim:** If $\xi_i > 0$ & $\xi_i^* > 0$ we get a contradiction

**Proof:** ① $\xi_i \geq \underbrace{y_i - \omega^T \phi(x_i) - b}_{k_i} - \epsilon$  $\left. \begin{array}{l} \min\limits_{\omega, \xi_i, \xi_i^*} c \sum (\xi_i + \xi_i^-) + \cdots \end{array} \right.$

$\xi_i = \max(0, k_i - \epsilon)$

② $\xi_i^* \geq \underbrace{\omega^T \phi(x_i) + b - y_i}_{-k_i} - \epsilon$

$\xi_i^* = \max(0, -k_i - \epsilon)$

If $\xi_i > 0$ then $\xi_i = k_i - \epsilon > 0 \Rightarrow k_i > \epsilon$  $\left. \begin{array}{l} \text{Contradiction} \end{array} \right.$

**Additionally**

If $\xi_i^* > 0$ then $\xi_i^* = -k_i - \epsilon > 0 \Rightarrow -k_i > \epsilon$

# SVR objective

- 1-norm Error, and $L_2$ regularized:

$$C \sum_i \left( \xi_i + \xi_i^* \right) + \frac{1}{2} \|w\|^2$$

$\frac{1}{2}$ is only to simplify future derivations & can be ignored

Instead of $C$ in the error, you could we a $\lambda$ in regularizer!

$C$ is inversely related to $\lambda$

s.t Constraints discussed earlier hold

# 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. $\forall i$,
      $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon + \xi_i$,
      $b + \mathbf{w}^\top \phi(\mathbf{x}_i) - y_i \leq \epsilon + \xi_i^*$,
      $\xi_i, \xi_i^* \geq 0$

- 2-norm Error, and $L_2$ regularized:

$$\sum_i \left( \xi_i + \xi_i^* \right)^2$$

# 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. $\forall i$,
    $y_i - \mathbf{w}^\top\phi(\mathbf{x}_i) - b \le \epsilon + \xi_i$,
    $b + \mathbf{w}^\top\phi(\mathbf{x}_i) - y_i \le \epsilon + \xi_i^*$,
    $\xi_i, \xi_i^* \ge 0$

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

*(handwritten annotation):* if $\xi_i < 0$ satisfies constraints so will $\xi_i = 0$ — since objective will be happier with smaller $\xi_i^2$, $\xi_i = 0$

# 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 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

# Regression through the eyes of Optimization

- **Unconstrained (Penalized) Optimization:** (Eg: Ridge)

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

- **Constrained Optimization 1:**

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

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

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

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

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

- **Equivalence**: $\lambda$ (Penalized) $\equiv \theta$ (Constrained)
- **Iteratively Solving:** Lasso, Regression with $L_0$ norm, Support Vector Regression
- **Duality**: Dual of Support Vector Regression $\Big\}$ Kernelization

- 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 :

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

where c is a constant.

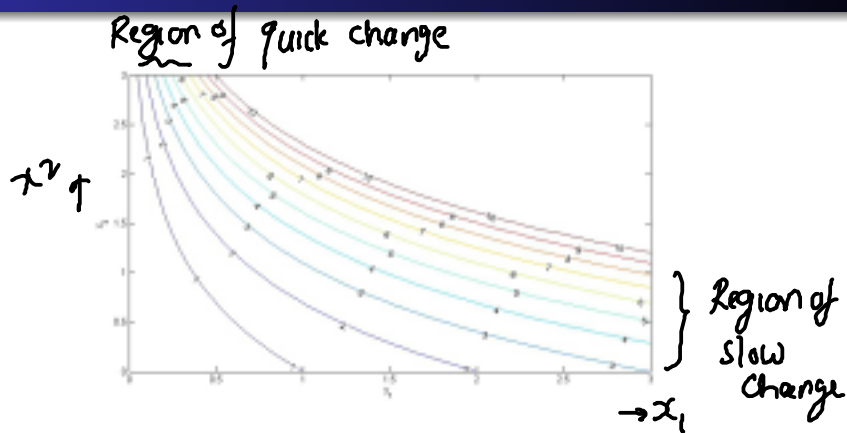Figure 1: 10 level curves for the function $f(x_1, x_2) = x_1 e^{x_2}$ (Figure 4.12 from https://www.cse.iitb.ac.in/~CS725/notes/classNotes/BasicsOfConvexOptimization.pdf)

- 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 :

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

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

Normalization

# Foundations: Gradient Vector

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

$$D_v(f, x) = V^T \nabla f(x)$$

$$\|\nabla f(x)\| = \max_V D_v(f, v)$$

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

Directional derivative along $x_i$

- <span style="color:red">Magnitude (euclidean norm)</span> of gradient vector at any point indicates maximum value of directional derivative at that point
- <span style="color:blue">Direction</span> of gradient vector indicates direction of this maximal directional derivative at that point. $\nabla f(x) / \|\nabla f(x)\|$
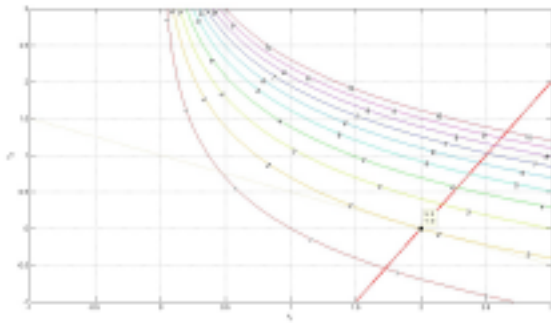
# Foundations: Gradient Vector



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

(Reminiscent of electrostatic flux)

# 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 ridge regression), ....

$$\mathcal{E}xpect: \nabla f\left(\omega_{MLE}\right) = 0$$

# 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 ridge regression), ....
- Ridge Regression: Find $\mathbf{w}$ such that

$$\mathbf{w}^* = \underset{\mathbf{w}}{\arg\min} \|\Phi\mathbf{w} - \mathbf{y}\|^2 + \lambda\|\mathbf{w}\|^2 \quad (5)$$

$$= \underset{\mathbf{w}}{\arg\min}(\mathbf{w}^T\Phi^T\Phi\mathbf{w} - 2\mathbf{w}^T\phi\mathbf{y} - \mathbf{y}^T\mathbf{y} + \lambda\|\mathbf{w}\|^2) \quad (6)$$

$O(\omega)$

$\nabla O(\omega^*) = 0$

- 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

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

$\implies \ldots\ldots\ldots$

- We would have

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

# Foundations: Necessary condition 1 (Solving Ridge)

- *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

$$f(\mathbf{w}) = (\mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\mathbf{w}^T \Phi^T \mathbf{y} - \mathbf{y}^T \mathbf{y} + \lambda ||\mathbf{w}||^2) \qquad (7)$$

$$\implies \nabla f(\mathbf{w}) = 2\Phi^T \Phi \mathbf{w} - 2\Phi^T \mathbf{y} + 2\lambda \mathbf{w} \qquad (8)$$

- We would have

$$\nabla f(\mathbf{w}^*) = 0 \qquad (9)$$

$$\implies 2(\Phi^T \Phi + \lambda I)\mathbf{w}^* - 2\Phi^T \mathbf{y} = 0 \qquad (10)$$

$$\implies \mathbf{w}^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{y} \qquad (11)$$