

Lecture 12: Support Vector Regression, Kernel Trick and Optimization Algorithm

Instructor: Prof. Ganesh Ramakrishnan

Some observations

- $\alpha_i, \alpha_i^* \geq 0, \mu_i, \mu_i^* \geq 0, \alpha_i + \mu_i = C$ and $\alpha_i^* + \mu_i^* = C$

Thus, $\alpha_i, \mu_i, \alpha_i^*, \mu_i^* \in [0, C], \forall i$

- If $0 < \alpha_i < C$, then $0 < \mu_i < C \rightarrow$ Support Vectors
(as $\alpha_i + \mu_i = C$)

- $\mu_i \xi_i = 0$ and $\alpha_i (y_i - w^T \phi(x_i) - b - \epsilon - \xi_i) = 0$ are complementary slackness conditions

↓ lying on ϵ -bdndry

So $0 < \alpha_i < C \Rightarrow \xi_i = 0$ and $y_i - w^T \phi(x_i) - b = \epsilon + \xi_i = \epsilon$

- ▶ All such points lie on the boundary of the ϵ band
- ▶ Using any point x_j (that is with $\alpha_j \in (0, C)$) on margin, we can recover b as:

$$b = y_j - w^T \phi(x_j) - \epsilon$$

} Assuming we know α_j & d_j

Support Vector Regression

Dual Objective

Dual function

- Let $L^*(\alpha, \alpha^*, \mu, \mu^*) = \min_{w, b, \xi, \xi^*} L(w, b, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*)$
- By weak duality theorem, we have:

$$\begin{aligned} \min_{w, b, \xi, \xi^*} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*) &\geq L^*(\alpha, \alpha^*, \mu, \mu^*) \\ \text{s.t. } y_i - w^\top \phi(x_i) - b &\leq \epsilon - \xi_i, \text{ and} \\ w^\top \phi(x_i) + b - y_i &\leq \epsilon - \xi_i^*, \text{ and} \\ \xi_i, \xi_i^* &\geq 0, \forall i = 1, \dots, n \end{aligned}$$

→ LHS is independent of $\alpha, \alpha^*, \mu, \mu^*$

- The above is true for any $\alpha_i, \alpha_i^* \geq 0$ and $\mu_i, \mu_i^* \geq 0$
- Thus,

$$\min_{w, b, \xi, \xi^*} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*) \geq \max_{\alpha, \alpha^*, \mu, \mu^*} L^*(\alpha, \alpha^*, \mu, \mu^*)$$

$$\begin{aligned} \text{s.t. } y_i - w^\top \phi(x_i) - b &\leq \epsilon - \xi_i, \text{ and} \\ w^\top \phi(x_i) + b - y_i &\leq \epsilon - \xi_i^*, \text{ and} \\ \xi_i, \xi_i^* &\geq 0, \forall i = 1, \dots, n \end{aligned}$$

Equality holds at KKT conditions under convexity

Dual objective

- In case of Support Vector Regression, we have a strictly convex objective and linear constraints \Rightarrow KKT conditions are necessary and sufficient and strong duality holds:

$$\min_{w, b, \xi, \xi^*} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*) = \max_{\alpha, \alpha^*, \mu, \mu^*} L^*(\alpha, \alpha^*, \mu, \mu^*)$$

s.t. $y_i - w^\top \phi(x_i) - b \leq \epsilon - \xi_i$, and
 $w^\top \phi(x_i) + b - y_i \leq \epsilon - \xi_i^*$, and
 $\xi_i, \xi_i^* \geq 0, \forall i = 1, \dots, n$

- This value is precisely obtained at the $(w, b, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*)$ that satisfies the necessary (and sufficient) optimality conditions
- Given strong duality, we can equivalently solve **KKT**

$$\max_{\alpha, \alpha^*, \mu, \mu^*} L^*(\alpha, \alpha^*, \mu, \mu^*)$$

- $$L(\alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\underline{\xi}_i + \underline{\xi}_i^*) + \sum_{i=1}^n (\alpha_i (\underline{y}_i - w^\top \phi(x_i) - \underline{b} - \underline{\epsilon} - \underline{\xi}_i) + \alpha_i^* (w^\top \phi(x_i) + \underline{b} - \underline{y}_i - \underline{\epsilon} - \underline{\xi}_i^*)) + \sum_{i=1}^n (\underline{\mu}_i \underline{\xi}_i + \underline{\mu}_i^* \underline{\xi}_i^*)$$

- We obtain w, b, ξ_i, ξ_i^* in terms of α, α^*, μ and μ^* by using the KKT conditions derived earlier as $w = \sum_{i=1}^n (\alpha_i - \alpha_i^*) \phi(x_i)$ and

$$\sum_{i=1}^n (\alpha_i - \alpha_i^*) = 0 \text{ and } \underline{\alpha}_i + \underline{\mu}_i = C \text{ and } \underline{\alpha}_i^* + \underline{\mu}_i^* = C$$

- Thus, we get:

$$\begin{aligned}
 & L(w, b, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*) \\
 &= \frac{1}{2} \sum_i \sum_j (\alpha_i - \alpha_i^*) (\alpha_j - \alpha_j^*) \phi^\top(x_i) \phi(x_j) + \sum_i (\underline{\xi}_i (C - \alpha_i - \mu_i) + \underline{\xi}_i^* (C - \alpha_i^* - \mu_i^*)) - \underline{b} \sum_i (\alpha_i - \alpha_i^*) - \underline{\epsilon} \sum_i (\alpha_i + \alpha_i^*) + \sum_i \underline{y}_i (\alpha_i - \alpha_i^*) - \sum_i \sum_j (\alpha_i - \alpha_i^*) (\alpha_j - \alpha_j^*) \phi^\top(x_i) \phi(x_j) \\
 &= -\frac{1}{2} \sum_i \sum_j (\alpha_i - \alpha_i^*) (\alpha_j - \alpha_j^*) \phi^\top(x_i) \phi(x_j) - \epsilon \sum_i (\alpha_i + \alpha_i^*) + \sum_i \underline{y}_i (\alpha_i - \alpha_i^*)
 \end{aligned}$$

Kernel function: $K(x_i, x_j) = \phi^T(x_i)\phi(x_j)$

- $w = \sum_{i=1}^n (\alpha_i - \alpha_i^*) \phi(x_i) \Rightarrow$ the final decision function
 $f(x) = w^T \phi(x) + b =$
 $\sum_{i=1}^n (\alpha_i - \alpha_i^*) \phi^T(x_i) \phi(x) + y_j - \sum_{i=1}^n (\alpha_i - \alpha_i^*) \phi^T(x_i) \phi(x_j) - \epsilon$
 x_j is any point with $\alpha_j \in (0, C)$
- The dual optimization problem to compute the α 's for SVR is:

$$\begin{aligned} \max_{\alpha_i, \alpha_i^*} & -\frac{1}{2} \sum_i \sum_j (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \underbrace{\phi^T(x_i)\phi(x_j)} \\ & -\epsilon \sum_i (\alpha_i + \alpha_i^*) + \sum_i y_i (\alpha_i - \alpha_i^*) \end{aligned}$$

s.t.

- ▶ $\sum_i (\alpha_i - \alpha_i^*) = 0$
- ▶ $\alpha_i, \alpha_i^* \in [0, C]$
- **We notice that the only way these three expressions involve ϕ is through $\phi^T(x_i)\phi(x_j) = K(x_i, x_j)$, for some i, j**

Kernelized form for SVR

- The *kernelized* dual optimization problem to compute the α 's for SVR is:

$$\begin{aligned} \max_{\alpha_i, \alpha_i^*} & -\frac{1}{2} \sum_i \sum_j (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) K(x_i, x_j) \\ & - \epsilon \sum_i (\alpha_i + \alpha_i^*) + \sum_i y_i (\alpha_i - \alpha_i^*) \end{aligned}$$

s.t.

- $\sum_i (\alpha_i - \alpha_i^*) = 0$
- $\alpha_i, \alpha_i^* \in [0, C]$

Q: What abt computing $f(x)$. Can we avoid ϕ & use only K in $f(x) = \omega^T \phi(x) + b$?

The Kernel function in SVR

- Again, invoking the **kernel function**:

$$K(x_1, x_2) = \phi^\top(x_1)\phi(x_2)$$

(eg: $(1 + x_i^\top x_j)^d$, $e^{-\frac{1}{2\sigma} \|x_i - x_j\|^2}$)

- The decision function becomes:

$$f(x) = \sum_i (\alpha_i - \alpha_i^*) K(x_i, x) + b$$

Polynomial
(finite ϕ)

(RBF) Radial Basis Function
(infinite ϕ)

- Using any point x_j (that is with $\alpha_j \in (0, C)$) on margin, we can recover b as:

$$b = y_j - w^\top \phi(x_j) - \epsilon = y_j - \sum_i (\alpha_i - \alpha_i^*) K(x_i, x_j)$$

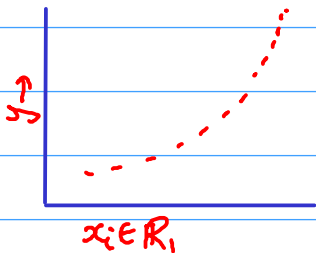
- Thus, the optimization problem as well as the final decision function are only in terms of the kernel function $K(x, x')$.
- We will see that, often, computing $K(x_1, x_2)$ does not even require computing $\phi(x_1)$ or $\phi(x_2)$ explicitly

Q1: Can ridge regression also be kernelized?

Q2: What are "valid" $K(\cdot, \cdot)$? Eg: $(1 + x_i^\top x_j)^d$

Idea behind Kernelization

Point $x_i \in \mathbb{R}^n$, whereas I want to predict in a high dimension $\phi(x_i)$ without actually enumerating ϕ and then computing $\phi^T(x_i)\phi(x_j)$



Can some $K(\cdot, \cdot)$ implicitly capture my non-linearity?

Possibly: $(1 + x_i x_j)^2 = K(x_i, x_j)$

How about Ridge Regression?

- Recall for Ridge Regression: $w = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$, where,

$$\phi(x_i) = \begin{bmatrix} \phi_1(x_i) \\ \phi_2(x_i) \\ \vdots \\ \phi_p(x_i) \end{bmatrix}$$

$$\Phi = \begin{bmatrix} \phi_1(x_1) & \dots & \phi_p(x_1) \\ \dots & \dots & \dots \\ \phi_1(x_m) & \dots & \phi_p(x_m) \end{bmatrix}$$

jth column was for $\phi_j(\cdot)$
ith row was for $\phi(x_i)$

and

$$(\Phi^T \Phi)_{(p,q)} = \sum_k \phi_p(x_k) \phi_q(x_k) \neq \sum_k \phi_k(x_p) \phi_k(x_q) = (\Phi \Phi^T)_{(p,q)}$$
$$y = \begin{bmatrix} y_1 \\ \dots \\ y_m \end{bmatrix}$$

- $(\Phi^T \Phi)_{ij} = \sum_{k=1}^m \phi_i(x_k) \phi_j(x_k)$ whereas
 $(\Phi \Phi^T)_{ij} = \sum_{k=1}^p \phi_k(x_i) \phi_k(x_j) = K(x_i, x_j)$

How about Ridge Regression?

- Given $w = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$ and using the identity
 $(P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} = P B^T (B P B^T + R)^{-1}$ (How abt $P=R=I$
 $\leftarrow B=\Phi$)
 - ▶ \Rightarrow
 - ▶ \Rightarrow

Scalar analog: $(\frac{1}{p} + b \frac{1}{r} b)^{-1} b \frac{1}{r} = p b (b p b + r)^{-1}$

How about Ridge Regression?

[Kernel Ridge Regression]

- Given $w = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$ and using the identity $(P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} = P B^T (B P B^T + R)^{-1}$

▶ $\Rightarrow w = \Phi^T (\Phi \Phi^T + \lambda I)^{-1} y = \sum_{i=1}^m \alpha_i \phi(x_i)$ where

$\alpha_i = \left((\Phi \Phi^T + \lambda I)^{-1} y \right)_i$ \rightarrow vector of $[\alpha_i]$

▶ \Rightarrow the final decision function

$f(x) = \phi^T(x) w = \sum_{i=1}^m \alpha_i \phi^T(x) \phi(x_i)$ $\rightarrow k(x, x_i)$

- Again, **We notice that the only way the decision function $f(x)$ involves ϕ is through $\phi^T(x_i) \phi(x_j)$, for some i, j**

$$\left. \begin{array}{l} R = I \\ P = I \\ B = \Phi \end{array} \right\}$$

The Kernel function in Ridge Regression

- We call $\phi^\top(x_1)\phi(x_2)$ a **kernel function**:
 $K(x_1, x_2) = \phi^\top(x_1)\phi(x_2)$
- The preceding expression for decision function becomes

$$f(x) = \sum_{i=1}^m \alpha_i K(x, x_i)$$

where $\alpha_i = (([K(x_i, x_j)] + \lambda I)^{-1} y)_i$

kernel ridge regression

Back to the Kernelized version of SVR

- The kernelized dual problem:

$$\begin{aligned} \max_{\alpha_i, \alpha_i^*} & -\frac{1}{2} \sum_i \sum_j (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \underline{K(x_i, x_j)} \\ & -\epsilon \sum_i (\alpha_i + \alpha_i^*) + \sum_i y_i (\alpha_i - \alpha_i^*) \end{aligned}$$

s.t.

- $\sum_i (\alpha_i - \alpha_i^*) = 0$
- $\alpha_i, \alpha_i^* \in [0, C]$

- The kernelized decision function:

$$f(x) = \sum_i (\alpha_i - \alpha_i^*) \underline{K(x_j, x)} + b$$

- Using any point x_j with $\alpha_j \in (0, C)$:

$$b = y_j - \sum_i (\alpha_i - \alpha_i^*) \underline{K(x_j, x_j)}$$

- Computing $K(x_1, x_2)$ often does not even require computing $\phi(x_1)$ or $\phi(x_2)$ explicitly

Main remaining question:
How to find α_i & α_i^* :

- ① Gradient ascent?
- ② Steepest ascent such as coordinate ascent on α_i / α_i^*
- ③ Constrained coordinate ascent?

Can I reduce α_i & α_i^* to β_i

An example

- Let $K(x_1, x_2) = (1 + x_1^\top x_2)^2$
- What $\phi(x)$ will give $\phi^\top(x_1)\phi(x_2) = K(x_1, x_2) = (1 + x_1^\top x_2)^2$
- Is such a ϕ guaranteed to exist?
- Is there a unique ϕ for given K ?

$$\phi(x) = ?$$

$$\begin{aligned} &= \left(1 + \sum_i x_{1i} x_{2i}\right)^2 \\ &= \left(1 + 2x_{11}x_{21} + 2x_{12}x_{22} \right. \\ &\quad \left. + 2x_{11}x_{21}x_{12} \right. \\ &\quad \left. x_{11}^2 x_{21}^2 + x_{12}^2 x_{22}^2\right) x_{22} \end{aligned}$$

- We can prove that such a ϕ exists
- For example, for a 2-dimensional x_i :

$$\phi(x_i) = \begin{bmatrix} 1 \\ x_{i1} \sqrt{2} \\ x_{i2} \sqrt{2} \\ x_{i1} x_{i2} \sqrt{2} \\ x_{i1}^2 \\ x_{i2}^2 \end{bmatrix}$$

$$\left. \begin{array}{l} \phi^\top(x_i) \phi(x_j) \\ = 1 + 2x_{i1}x_{j1} + 2x_{i2}x_{j2} \\ + 2x_{i1}x_{i2}x_{j1}x_{j2} \\ + x_{i1}^2x_{j1}^2 + x_{i2}^2x_{j2}^2 \end{array} \right\}$$

- $\phi(x_i)$ exists in a 5-dimensional space
- Thus, to compute $K(x_1, x_2)$, all we need is $x_1^\top x_2$, and there is no need to compute $\phi(x_i)$

H/W: If $K(x_i, x_j) = (1 + x_i x_j)^d$ what would the required dimension of ϕ be?

Introduction to the Kernel Trick (more later)

- **Kernels** operate in a *high-dimensional, implicit* feature space without ever computing the coordinates of the data in that space, but rather by simply computing the Kernel function
- This approach is called the "*kernel trick*" and will talk about *valid kernels* a little later...
- This operation is often computationally cheaper than the explicit computation of the coordinates

Solving the SVR Dual Optimization Problem

- The SVR dual objective is:

$$\max_{\alpha_i, \alpha_i^*} -\frac{1}{2} \sum_i \sum_j (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) K(x_i, x_j) \\ - \epsilon \sum_i (\alpha_i + \alpha_i^*) + \sum_i y_i (\alpha_i - \alpha_i^*)$$

- This is a linearly constrained quadratic program (LCQP), just like the constrained version of Lasso
- There exists no closed form solution to this formulation
- Standard QP (LCQP) solvers¹ can be used
- Question: Are there more specific and efficient algorithms for solving SVR in this form?

¹https://en.wikipedia.org/wiki/Quadratic_programming#Solvers_and_scripting_.28programming.29_languages

Solving the SVR Dual Optimization Problem

- It can be shown that the objective:

$$\max_{\alpha_i, \alpha_i^*} -\frac{1}{2} \sum_i \sum_j (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) K(x_i, x_j) - \epsilon \sum_i (\alpha_i + \alpha_i^*) + \sum_i y_i (\alpha_i - \alpha_i^*)$$

- can be written as:

$$\max_{\beta_i} -\frac{1}{2} \sum_i \sum_j \beta_i \beta_j K(x_i, x_j) - \epsilon \sum_i |\beta_i| + \sum_i y_i \beta_i$$

s.t.

- ▶ $\sum_i \beta_i = 0$
- ▶ $\beta_i \in [-C, C], \forall i$

- Even for this form, standard QP (LCQP) solvers² can be used
- Question: How about (iteratively) solving for two β_i 's at a time?

(Since $\sum \beta_i = 0$, coordinate ascent on a β_i will not work)

- ▶ This is the idea of the Sequential Minimal Optimization (SMO) algorithm

²https://en.wikipedia.org/wiki/Quadratic_programming#Solvers_and_scripting_.28programming.29_languages

Sequential Minimal Optimization (SMO) for SVR

- Consider:

$$\max_{\beta_i} -\frac{1}{2} \sum_i \sum_j \beta_i \beta_j K(x_i, x_j) - \epsilon \sum_i |\beta_i| + \sum_i y_i \beta_i$$

s.t.

- $\sum_i \beta_i = 0$
- $\beta_i \in [-C, C], \forall i$

- The SMO subroutine can be defined as:

- 1 Initialise $\beta_1, \dots, \beta_{n-1}$ to some value $\in [-C, C]$ & let $\beta_n = -\sum_{i=1}^{n-1} \beta_i$
- 2 Pick β_i, β_j to estimate closed form expression for next iterate (i.e. $\beta_i^{new}, \beta_j^{new}$)
- 3 Check if the KKT conditions are satisfied

- ★ If not, choose β_i and β_j that worst violate the KKT conditions and reiterate

(see tut 5) $\beta_i^{new} + \beta_j^{new} = \beta_i + \beta_j \Rightarrow \beta_j^{new} = \beta_i + \beta_j - \beta_i^{new}$

In the objective substitute for β_j^{new} to remove constraint

