Introduction to Machine Learning - CS725
Instructor: Prof. Ganesh Ramakrishnan
Lecture 12 - Support Vector Regression and its
Dual using Optimization Principles, Kernel Trick

Recap: Lagrange Function for SVR

- $\min_{\mathbf{w},b,\xi_{i},\xi_{i}^{*}} \frac{1}{2} \|\mathbf{w}\|^{2} + C \sum_{i} (\xi_{i} + \xi_{i}^{*}) \\
 \text{s.t. } \forall i, \\
 y_{i} \mathbf{w}^{\top} \phi(\mathbf{x}_{i}) b \leq \epsilon + \xi_{i}, \\
 b + \mathbf{w}^{\top} \phi(\mathbf{x}_{i}) y_{i} \leq \epsilon + \xi_{i}^{*}, \\
 \xi_{i}, \xi_{i}^{*} > 0$
- Consider corresponding lagrange multipliers α_i , α_i^* , μ_i and μ_i^*
- The Lagrange Function is $L(\mathbf{w}, \alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i} (\xi_i + \xi_i^*) + \sum_{i=1}^{m} \alpha_i (y_i \mathbf{w}^\top \phi(\mathbf{x}_i) b \epsilon \xi_i) + \sum_{i=1}^{m} \alpha_i^* (b + \mathbf{w}^\top \phi(\mathbf{x}_i) y_i \epsilon \xi_i^*) \sum_{i=1}^{m} \mu_i \xi_i \sum_{i=1}^{m} \mu_i^* \xi_i^*$



```
Recap: KKT conditions for the Constrained (Convex) Problem Assume the on values of \left\{\hat{\mathbf{w}},\hat{b},\hat{\xi},\hat{\xi}^*,\hat{\alpha},\hat{\alpha}^*,\hat{\mu},\hat{\mu}^*\right\} at KKT when not explicitly specified
```

Recap: Necessary and Sufficient SVR KKT conditions

- Differentiating the Lagrangian w.r.t. \mathbf{w} , $\mathbf{w} \alpha_i \phi(\mathbf{x}_i) + \alpha_i^* \phi(\mathbf{x}_i) = 0$ i.e. $\mathbf{w} = \sum_{i=1}^m (\alpha_i \alpha_i^*) \phi(\mathbf{x}_i)$
- Differentiating the Lagrangian w.r.t. ξ_i , $C \alpha_i \mu_i = 0$ i.e. $\alpha_i + \mu_i = C$
- Differentiating the Lagrangian w.r.t ξ_i^* , $\alpha_i^* + \mu_i^* = C$
- Differentiating the Lagrangian w.r.t b, $\sum_{i}^{m} (\alpha_{i}^{*} \alpha_{i}) = 0$
- Complimentary slackness: $\alpha_i(\mathbf{y}_i \mathbf{w}^{\top} \phi(\mathbf{x}_i) b \epsilon \xi_i) = 0$





For any point (x_i, y_i) , the product $\alpha_i \alpha_i^* = 0$. Suppose $a_i > 0$ $a_i > 0$ + $a_i = a_i + a_i = 0$ Suppose $a_i > 0$ $a_i > 0$ By complementary $a_i = a_i + a_i = 0$ Slackness

Slackness - 26 = 8: + 8: : KKT includes original constraint set as well E: 18: 20 A contradiction!

- For any point (\mathbf{x}_i, y_i) , the product $\alpha_i \alpha_i^* = 0$.
 - Let $\alpha_i > 0$ and $\alpha_i^* > 0$. This leads to a contradiction.
 - By Complimentary slackness, $y_i \mathbf{w}^{\top} \phi(\mathbf{x}_i) b \epsilon \xi_i = 0$ AND $b + \mathbf{w}^{\top} \phi(\mathbf{x}_i) y_i \epsilon \xi_i^* = 0$. Adding up the two equalities gives us: $\xi_i + \xi_i^* = -2\epsilon$.
 - Since only one of ξ_i and ξ_i^* can be non-zero, \Longrightarrow the non-zero component is negative, which is a contradiction since $\xi_i, \xi_i^* \geq 0$
 - Thus, $\alpha_i \alpha_i^* \propto \max\{\alpha_i, \alpha_i^*\}$
- For points within the ϵ -insensitive tube $\alpha_i = 0$ and $\alpha_i^* = 0$:

We saw this last time



- For any point (x_i, y_i) , the product $\alpha_i \alpha_i^* = 0$.
 - Let $\alpha_i > 0$ and $\alpha_i^* > 0$. This leads to a contradiction.
 - By Complimentary slackness, $y_i \mathbf{w}^{\top} \phi(\mathbf{x}_i) b \epsilon \xi_i = 0$ AND $b + \mathbf{w}^{\top} \phi(\mathbf{x}_i) y_i \epsilon \xi_i^* = 0$. Adding up the two equalities gives us: $\xi_i + \xi_i^* = -2\epsilon$.
 - Since only one of ξ_i and ξ_i^* can be non-zero, \Longrightarrow the non-zero component is negative, which is a contradiction since $\xi_i, \xi_i^* \geq 0$
 - Thus, $\alpha_i \alpha_i^* \propto \max\{\alpha_i, \alpha_i^*\}$
- For points within the ϵ -insensitive tube $\alpha_i = 0$ and $\alpha_i^* = 0$:
 - If $y_i \mathbf{w}^{\top} \phi(\mathbf{x}_i) b \epsilon \xi_i < 0$, then $\alpha_i = 0$, $\mu_i = C$ and $\xi_i = 0$. Similarly, $b + \mathbf{w}^{\top} \phi(\mathbf{x}_i) - y_i - \epsilon < 0$ leading to $\alpha_i^* = 0$.

• $\alpha_i = C$ and $\alpha_i^* = C$ correspond to points lying either outside or on the ϵ -tube:

We have seen this

- $\alpha_i = C$ and $\alpha_i^* = C$ correspond to points lying either outside or on the ϵ -tube:
 - If $\alpha_i = C$, then $\mu_i = 0$ and $y_i \mathbf{w}^{\top} \phi(\mathbf{x}_i) b \epsilon = \xi_i \geq 0$.
 - Similarly, $\alpha_i^* = C$ corresponds to points lying below (or beyond) the lower ϵ -band.
- For points on boundary of the ϵ -insensitive tube $\alpha_i \in [0, C]$:

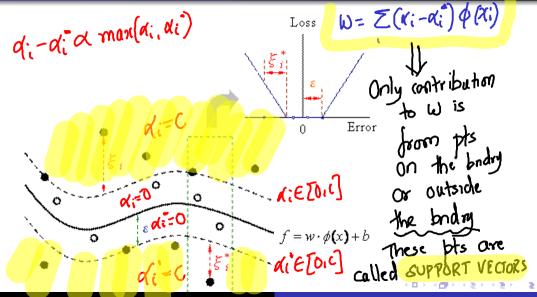
We argued that
$$\alpha; \epsilon(0,C)$$
 should correspond to pts "on the boundary" of ϵ -tube!



- $\alpha_i = C$ and $\alpha_i^* = C$ correspond to points lying either outside or on the ϵ -tube:
 - If $\alpha_i = C$, then $\mu_i = 0$ and $y_i \mathbf{w}^\top \phi(\mathbf{x}_i) b \epsilon = \xi_i \ge 0$.
 - Similarly, $\alpha_i^* = C$ corresponds to points lying below (or beyond) the lower ϵ -band.
- For points on boundary of the ϵ -insensitive tube $\alpha_i \in [0, C]$:
 - For any point on the upper margin, $y_i \mathbf{w}^{\top} \phi(\mathbf{x}_i) b \epsilon = 0$ and $\xi_i = 0 \Longrightarrow \mu_i \ge 0 \Longrightarrow \alpha_i \in [0, C]$. Similarly, $\alpha_i^* \in [0, C]$ for points lying on the margin of the lower ϵ -band.



Support Vector Regression (SVR)



Recap: Retrieving solution for b

- $\mu_i \xi_i = 0$ and $\alpha_i (y_i \mathbf{w}^{\top} \phi(\mathbf{x}_i) b \epsilon \xi_i) = 0$ are complementary slackness conditions So $0 < \alpha_i < C \Rightarrow \xi_i = 0$ and $y_i - \mathbf{w}^{\top} \phi(\mathbf{x}_i) - b = \epsilon + \xi_i = \epsilon$
 - \bullet All such points lie on the boundary of the ϵ band
 - Using any point \mathbf{x}_j (that is with $\alpha_j \in (0, C)$) on margin, we can recover b as:

$$b = y_j - \mathbf{w}^{\top} \phi(\mathbf{x}_j) - \epsilon$$

Hereafter we will not bother abt b much



Support Vector Regression Dual Objective

Weak Duality and SVR

- Defined for $A,A,A,A^* \ge 0$ $L^*(\alpha,\alpha^*,\mu,\mu^*) = \min_{\mathbf{w},b,\xi,\xi^*} L(\mathbf{w},b,\xi,\xi^*,\alpha,\alpha^*,\mu,\mu^*)$
- By weak duality theorem, for any $\alpha_i, \alpha_i^* \geq 0$ and $\mu_i, \mu_i^* \geq 0$:

Weak Duality and SVR

- $L^*(\alpha, \alpha^*, \mu, \mu^*) = \min_{\mathbf{w}, \mathbf{b}, \xi, \xi^*} L(\mathbf{w}, \mathbf{b}, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*)$
- By weak duality theorem, for any $\alpha_i, \alpha_i^* \geq 0$ and $\mu_i, \mu_i^* \geq 0$: $\min_{\mathbf{w}, b, \xi, \xi^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) \geq L^*(\alpha, \alpha^*, \mu, \mu^*) \text{ s.t. } \mathbf{a}_i \mathbf{a}_i$
- Thus,



Weak Duality and SVR

- $L^*(\alpha, \alpha^*, \mu, \mu^*) = \min_{\mathbf{w}, b, \xi, \xi^*} L(\mathbf{w}, b, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*)$
- By weak duality theorem, for any $\alpha_i, \alpha_i^* \geq 0$ and $\mu_i, \mu_i^* \geq 0$: $\min_{\mathbf{w}, b, \xi, \xi^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) \geq L^*(\alpha, \alpha^*, \mu, \mu^*)$ s.t. $y_i \mathbf{w}^\top \phi(\mathbf{x}_i) b \leq \epsilon \xi_i$, and $\mathbf{w}^\top \phi(\mathbf{x}_i) + b y_i \leq \epsilon \xi_i^*$ and $\xi_i, \xi^* \geq 0$. $\forall i = 1, \dots, n$
- Thus, $\min_{\mathbf{w},b,\xi,\xi^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) \ge \max_{\alpha,\alpha^*,\mu,\mu^*} L^*(\alpha,\alpha^*,\mu,\mu^*)$
 - s.t. $y_i \mathbf{w}^{\top} \phi(\mathbf{x}_i) b \leq \epsilon \xi_i$, and $\mathbf{w}^{\top} \phi(\mathbf{x}_i) + b y_i \leq \epsilon \xi_i^*$ and $\xi_i, \xi^* > 0$, $\forall i = 1, ..., n$



SVR Dual objective

- Assume: By convexity, KKT conditions are necessary and sufficient and strong duality holds (for $\alpha, \alpha^* \geq 0$): Now $\min_{\mathbf{w}, b, \xi, \xi^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) = \max_{\alpha, \alpha^*, \mu, \mu^*} L^*(\alpha, \alpha^*, \mu, \mu^*)$ s.t. $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon - \xi_i$, and $w^\top \phi(\mathbf{x}_i) + b - y_i \leq \epsilon - \xi_i^*$ and $\xi_i, \xi^* > 0, \forall i = 1, \ldots, n$
- This value is precisely obtained at the $\left\{\hat{\mathbf{w}},\hat{b},\hat{\xi},\hat{\xi}^*,\hat{\alpha},\hat{\alpha}^*,\hat{\mu},\hat{\mu}^*\right\}$ that satisfies the necessary (and sufficient) KKT optimality conditions [KKT Constraint Set]



SVR Dual objective (contd)

• For $\alpha, \alpha^* \geq 0$ and $\left\{\hat{\mathbf{w}}, \hat{b}, \hat{\xi}, \hat{\xi}^*, \hat{\alpha}, \hat{\alpha}^*, \hat{\mu}, \hat{\mu}^*\right\}$ from [KKT Constraint Set]: $\min_{\mathbf{w}, b, \xi, \xi^*} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) = \max_{\alpha, \alpha^*, \mu, \mu^*} L^*(\alpha, \alpha^*, \mu, \mu^*)$ s.t. $y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b \leq \epsilon - \xi_i$, and $w^\top \phi(\mathbf{x}_i) + b - y_i \leq \epsilon - \xi_i^*$ and $\xi_i, \xi^* \geq 0$, $\forall i = 1, \dots, n$

• Given strong duality, we can equivalently solve: $\max_{\hat{\alpha},\hat{\alpha}^*,\hat{\mu},\hat{\mu}^*} L^*(\hat{\alpha},\hat{\alpha}^*,\hat{\mu},\hat{\mu}^*)$



• $L(\hat{\alpha}, \hat{\alpha}^*, \hat{\mu}, \hat{\mu}^*) = \frac{1}{2} \|\hat{\mathbf{w}}\|^2 + C \sum_{i=1}^m (\hat{\xi}_i + \hat{\xi}_i^*) + C \sum_{i=1}^m (\hat{\xi$ $\sum_{i=1}^{m} \left(\hat{\alpha}_{i} (\mathbf{y}_{i} - \mathbf{w}^{\top} \phi(\mathbf{x}_{i}) - b - \epsilon - \hat{\xi}_{i}) + \hat{\alpha}_{i}^{*} (\mathbf{w}^{\top} \phi(\mathbf{x}_{i}) + b - \mathbf{y}_{i} - \epsilon - \hat{\xi}_{i}^{*}) \right)$ $\sum_{i=1}^{m} (\hat{\mu}_i \hat{\xi}_i + \hat{\mu}_i^* \hat{\xi}_i^*)$ Substitut KKT into L(...)

• We obtain $\hat{\mathbf{w}}$, \hat{b} , $\hat{\xi}_i$, $\hat{\xi}_i^*$ in terms of $\hat{\alpha}$, $\hat{\alpha}^*$, $\hat{\mu}$ and $\hat{\mu}^*$ by using the KKT conditions derived earlier as $\hat{\mathbf{w}} = \sum_{i=1}^{m} (\hat{\alpha}_i - \hat{\alpha}_i^*) \phi(\mathbf{x}_i)$ and $\sum (\hat{lpha}_i - \hat{lpha}_i^*) = 0$ and $\hat{lpha}_i + \hat{\mu}_i = C$ and $\hat{lpha}_i^* + \hat{\mu}_i^* = C$

and
$$\sum (\hatlpha_i-\hatlpha_i^*)=$$
 0 and $\hatlpha_i+\hat\mu_i=$ ${\cal C}$ and $\hatlpha_i^*+\hat\mu_i^*={\cal C}$

Dropping the messy
$$\hat{}$$
 hat notation...

•
$$L(\alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\mathbf{x}_i + \mathbf{x}_i^*) + \sum_{i=1}^m (\alpha_i (\mathbf{y}_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - \mathbf{b} - \epsilon - \mathbf{x}_i) + \alpha_i^* (\mathbf{w}^\top \phi(\mathbf{x}_i) + \mathbf{b} - \mathbf{y}_i - \epsilon - \mathbf{x}_i^*)$$

$$\sum_{i=1}^{\infty} (lpha_i(\mathbf{y}_i - \mathbf{w}^*) \mathbf{q})$$

• Invoking
$$\mathbf{w} = \sum_{i=1}^{m} (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i)$$
 and $\sum_{i=1}^{m} (\alpha_i - \alpha_i^*) = 0$ and $\alpha_i + \mu_i = C$ and $\alpha_i^* + \mu_i^* = C$, we get simplify using

$$\alpha_{i} + \mu_{i} = C \text{ and } \alpha_{i}^{*} + \mu_{i}^{*} = C, \text{ we get simplify using}$$

$$\sum_{i} \mathcal{E}_{i} \left(-\alpha_{i}^{*} - \mu_{i}^{*} \right) = O \sum_{i} \frac{1}{2} \left[\mu_{i} \right]_{i}^{2} \frac{1}{2} \sum_{i} \left[\alpha_{i}^{*} - \alpha_{i}^{*} \right] \left[\alpha_{j}^{*} - \alpha_{j}^{*} \right]$$

Dropping the messy : hat notation...

•
$$L(\alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + C \sum_{i=1}^m (\xi_i + \xi_i^*$$

$$\sum_{i=1}^{m} \left(\alpha_{i}(y_{i} - \mathbf{w}^{\top}\phi(\mathbf{x}_{i}) - b - \epsilon - \xi_{i}) + \alpha_{i}^{*}(\mathbf{w}^{\top}\phi(\mathbf{x}_{i}) + b - y_{i} - \epsilon - \xi_{i}^{*})\right)$$

$$\sum_{i=1}^{m} \left(\mu_{i}\xi_{i} + \mu_{i}^{*}\xi_{i}^{*}\right)$$

• Invoking $\mathbf{w} = \sum_{i=1}^{m} (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i)$ and $\sum_{i=1}^{m} (\alpha_i - \alpha_i^*) = 0$ and

$$lpha_i + \mu_i = C$$
 and $lpha_i^* + \mu_i^* = C$, we get $L(\mathbf{w}, b, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \sum_i \sum_j (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_i^*)$

 $L(\mathbf{w}, b, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \sum_i \sum_i (\alpha_i - \alpha_i^*)(\alpha_i - \alpha_i^*)$ $(\alpha_i^*)\phi^{\top}(\mathbf{x}_i)\phi(\mathbf{x}_i) + \sum_i (\xi_i(C - \alpha_i - \mu_i) + \xi_i^*(C - \alpha_i^* - \mu_i^*)) - \xi_i^*(C - \alpha_i^* - \mu_i^*)$ $b\sum_{i}(\alpha_{i}-\alpha_{i}^{*})-\epsilon\sum_{i}(\alpha_{i}+\alpha_{i}^{*})+\sum_{i}y_{i}(\alpha_{i}-\alpha_{i}^{*}) \frac{1}{2}\sum_{i}\sum_{j}(\alpha_{i}-\alpha_{i}^{*})(\alpha_{j}-\alpha_{j}^{*})\phi^{\top}(\mathbf{x}_{i})\phi(\mathbf{x}_{j})$ [Only in terms of

Developing further...

•
$$L(\alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \sum_{i=1}^m (\alpha_i (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) - b - \epsilon - \xi_i) + \alpha_i^* (\mathbf{w}^\top \phi(\mathbf{x}_i) + b - y_i - \epsilon - \xi_i^*)$$

$$\sum_{i=1}^{m} (\mu_i \xi_i + \mu_i^* \xi_i^*)$$

•
$$L(\mathbf{w}, b, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \sum_i \sum_j (\alpha_i - \alpha_i^*) (\alpha_j - \alpha_j^*) \phi^\top(\mathbf{x}_i) \phi(\mathbf{x}_j) + \sum_i (\xi_i (C - \alpha_i - \mu_i) + \xi_i^* (C - \alpha_i^* - \mu_i^*)) - b \sum_i (\alpha_i - \alpha_i^*) - \epsilon \sum_i (\alpha_i + \alpha_i^*) + \sum_i y_i (\alpha_i - \alpha_i^*) - \sum_i \sum_i (\alpha_i - \alpha_i^*) (\alpha_i - \alpha_i^*) \phi^\top(\mathbf{x}_i) \phi(\mathbf{x}_i)$$

Developing further...

•
$$L(\mathbf{w}, b, \xi, \xi^*, \alpha, \alpha^*, \mu, \mu^*) = \frac{1}{2} \sum_i \sum_j (\alpha_i - \alpha_i^*) (\alpha_j - \alpha_j^*) \phi^\top(\mathbf{x}_i) \phi(\mathbf{x}_j) + \sum_i (\xi_i (C - \alpha_i - \mu_i) + \xi_i^* (C - \alpha_i^* - \mu_i^*)) - b \sum_i (\alpha_i - \alpha_i^*) - \epsilon \sum_i (\alpha_i + \alpha_i^*) + \sum_i y_i (\alpha_i - \alpha_i^*) - \sum_i \sum_j (\alpha_i - \alpha_i^*) (\alpha_j - \alpha_j^*) \phi^\top(\mathbf{x}_i) \phi(\mathbf{x}_j)$$

$$= -\frac{1}{2} \sum_i \sum_j (\alpha_i - \alpha_i^*) (\alpha_j - \alpha_j^*) \phi^\top(\mathbf{x}_i) \phi(\mathbf{x}_j) - \epsilon \sum_i (\alpha_i + \alpha_i^*) + \sum_i y_i (\alpha_i - \alpha_i^*)$$

SVR Dual using only dot products $\phi^T(\mathbf{x}_i)\phi(\mathbf{x}_j)$

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

SVR Dual using only dot products $\phi^T(\mathbf{x}_i)\phi(\mathbf{x}_j)$

- $\mathbf{w} = \sum_{i=1}^{m} (\alpha_i \alpha_i^*) \phi(\mathbf{x}_i) \Rightarrow$ the final decision function $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b = \sum_{i=1}^{m} (\alpha_i \alpha_i^*) \phi^T(\mathbf{x}_i) \phi(\mathbf{x}) + y_j \sum_{i=1}^{m} (\alpha_i \alpha_i^*) \phi^T(\mathbf{x}_i) \phi(\mathbf{x}_j) \epsilon \mathbf{x}_j$ is any point with $\alpha_j \in (0, C)$.
- The dual optimization problem to compute the α 's for SVR is:
 - $\max_{\alpha_i, \alpha_i^*} \frac{1}{2} \sum_i \sum_j (\alpha_i \alpha_i^*) (\alpha_j \alpha_j^*) \phi^{\top}(\mathbf{x}_i) \phi(\mathbf{x}_j) \epsilon \sum_i (\alpha_i + \alpha_i^*) + \sum_i y_i (\alpha_i \alpha_i^*)$ • $\mathbf{s.t} \sum_i (\alpha_i - \alpha_i^*) = 0 \& \alpha_i, \alpha_i^* \in [0, C]$
- We notice that the only way these three expressions involve ϕ is through $\phi^{\top}(\mathbf{x}_i)\phi(\mathbf{x}_j) = K(\mathbf{x}_i,\mathbf{x}_j)$, for some i,j



Kernel function: $K(\mathbf{x}_i, \mathbf{x}_j) = \phi^T(\mathbf{x}_i)\phi(\mathbf{x}_j)$

- We call $\phi^{\top}(\mathbf{x}_i)\phi(\mathbf{x}_j)$ a kernel function: $K(\mathbf{x}_i, \mathbf{x}_j) = \phi^{\top}(\mathbf{x}_i)\phi(\mathbf{x}_j)$
- The Kernel Trick: For some important choices of ϕ , compute $K(\mathbf{x}_i, \mathbf{x}_j)$ directly and more efficiently than having to explicitly compute/enumerate $\phi(\mathbf{x}_i)$ and $\phi(\mathbf{x}_j)$
- The expression for decision function becomes $f(x) = \sum_{i=1}^{m} \alpha_i K(\mathbf{x}, \mathbf{x}_i) \qquad \text{Similar 5 of query plants}$ Computation of a six an additional to the six and th
- Computation of α_i is specific to the objective function being minimized: Closed form exists for Ridge regression but NOT for SVR

 Ken nel Ridge Rogerssion (Tut 5)

The Kernelized version of SVR

• The kernelized dual problem:

$$\max_{\alpha_i,\alpha_i^*} -\frac{1}{2} \sum_i \sum_j (\alpha_i - \alpha_i^*) (\alpha_j - \alpha_j^*) \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$$
$$-\epsilon \sum_i (\alpha_i + \alpha_i^*) + \sum_i y_i (\alpha_i - \alpha_i^*)$$

- such that $\sum_{i}(\alpha_{i}-\alpha_{i}^{*})=0$ and $\alpha_{i},\alpha_{i}^{*}\in[0,C]$
- Kernelized decision function: $f(\mathbf{x}) = \sum_{i} (\alpha_i \alpha_i^*) K(\mathbf{x}_i, \mathbf{x}) + b$
- Using any \mathbf{x}_j with $\alpha_j \in (0, C)$: $b = y_j \sum_i (\alpha_i \alpha_i^*) K(\mathbf{x}_i, \mathbf{x}_j)$
- Computing $K(\mathbf{x}_1, \mathbf{x}_2)$ often does not even require computing $\phi(\mathbf{x}_1)$ or $\phi(\mathbf{x}_2)$ explicitly



Tutorial 5: Derive kernelized expression for Ridge Regression

Tutorial 5: Kernelizing Ridge Regression

$$\overline{\Phi} = \begin{bmatrix} \phi(x_i) \\ \phi(x_m) \end{bmatrix} \quad \overline{\Phi}^T \underline{\Phi} = \begin{bmatrix} \sum_k \phi_i(x_k) \phi_j(x_k) \end{bmatrix} \neq \begin{bmatrix} K(x_i, x_j) \\ \end{bmatrix}$$

- Given $w = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$ and using the identity $(P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} = PB^T (BPB^T + R)^{-1}$ (i,j) tentry
 - $\Rightarrow w = \Phi^T (\underline{\Phi}\underline{\Phi}^T + \lambda I)^{-1} y = \sum_{i=1}^m \alpha_i \phi(x_i)$ where $\alpha_i = ((\Phi \underline{\Phi}^T + \lambda I)^{-1} y)_i$
 - \Rightarrow the final decision function $f(\mathbf{x}) = \phi^T(\mathbf{x})\mathbf{w} = \sum_{i=1}^m \alpha_i \phi^T(\mathbf{x})\phi(\mathbf{x}_i)$
- Again, We notice that the only way the decision function $f(\mathbf{x})$ involves ϕ is through $\phi^{\top}(\mathbf{x}_i)\phi(\mathbf{x}_j)$, for some i,j

Hint: Try identity for P.B.R EIR (ic scalars)

$$\Phi \Phi^{T} = K(\pi_{i}, \pi_{j})$$

Basis function expansion and the Kernel trick

• We began with functional form called basis function expansion¹

$$f(\mathbf{x}) = \sum_{j=1}^{p} w_j \phi_j(\mathbf{x})$$

And landed up with an equivalent form for Ridge and SVR

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

• Aside: For $p \in [0, \infty)$, with what K, kind of regularizers, loss functions, *etc.*, will these dual representations hold?²



¹Each ϕ_j is called a *basis function*. *b* can be absorbed in ϕ . See Section 2.8.3 of Tibshi

²Section 5.8.1 of Tibshi.

The Representer Theorem & Reproducing Kernel Hilbert Space (RKHS) [Optional Slide]

① The solution $f^* \in \mathcal{H}_K$ (Hilbert space) to the following problem

$$\int_{\mathbf{mcreasing}}^{\mathbf{r}} f^* = \underset{f \in \mathcal{H}_{M_{\mathbf{K}}}}{\operatorname{arg\,min}} \sum_{i=1}^{m} \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|f\|_{K})$$

can be always written as $f^*(\mathbf{x}) = \sum_{i=1}^m \alpha_i K(\mathbf{x}, \mathbf{x}^{(i)})$, provided $(x,x) = x^{T}x = ||x||_{2}^{2} \quad \forall x \in \mathbb{R}^{n} \quad \text{ such edian space}$ $(x,x) = x^{T}x = ||x||_{2}^{2} \quad \forall x \in \mathbb{R}^{n} \quad \text{ such edian space}$ $(x,x) = x^{T}x = ||x||_{2}^{2} \quad \forall x \in \mathbb{R}^{n} \quad \text{ such edian space}$ $(x,x) = x^{T}x = ||x||_{2}^{2} \quad \forall x \in \mathbb{R}^{n} \quad \text{ such edian space}$

The Representer Theorem & Reproducing Kernel Hilbert Space (RKHS)

• More specifically, if $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$ and $K(\mathbf{x}',\mathbf{x}) = \phi^T(\mathbf{x})\phi(\mathbf{x}')$ then the solution $\mathbf{w}^* \in \Re^n$ to the following problem $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) + \Omega(\|\mathbf{w}\|_2)$ $\nabla^{\text{ont}}(\mathbf{w}^*, b^*) = \arg\min_{\mathbf{w}, b} \sum_{i=1}^m \mathbf{E}\left(f\left(\mathbf{x}^{(i)}\right), y^{(i)}\right) +$ $\|\mathbf{w}\|_{2}$. \Re^{n} is the Hilbert space and $K(.,\mathbf{x}):\mathcal{X}\to\Re$ is the Reproducing (RKHS) Kernel