

SVM & SVR

Support Vector Machine & Support Vector Regression

- Maximum margin classifier
- Primal and dual forms
- Linearly separable data
- Not linearly separable data
- Slack variables and hinge loss

- Epsilon - insensitive loss

When we use SVM:

Sparse data
(include under-determined systems)
 $N_c < (3-10) \times \text{d.o.f.}$

Adding constraints to learning problem

Non-linear mapping to expanded feature space

Use kernels $k(x, x')$

I. Concept (non-augmented notation & space)

① Non-linear mapping $\underline{x} \xrightarrow{\text{N.L. mapping}} \underline{u} = \phi(\underline{x})$

← Assume data is linearly separable

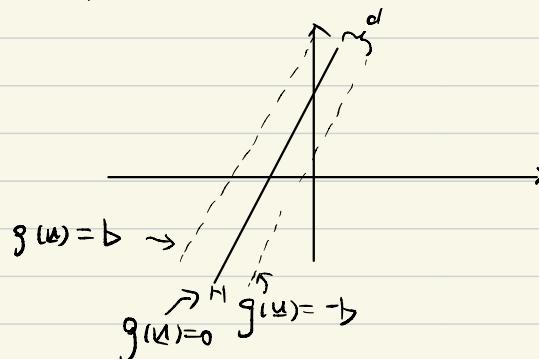
② Optimized version of using a margin.

$$g(\underline{x}) = \underline{w}^T \underline{x} + w_0 \geq b > 0 \text{ for } \underline{x} \in S_1$$

$$g'(\underline{u}) = \underline{w}^T \underline{u} + w_0 \geq b' > 0 \text{ for } \underline{u} \in S_1 \Leftrightarrow \phi(\underline{x}) \text{ or } k(\underline{x}, \underline{x})$$

unaugmented

$$\Rightarrow g(\underline{u}_n) = \begin{cases} \underline{w}^T \underline{u}_n + w_0 \geq b > 0 & \forall \underline{u}_n \in S_1 \\ \underline{w}^T \underline{u}_n + w_0 \leq -b < 0 & \forall \underline{u}_n \in S_2 \end{cases}$$

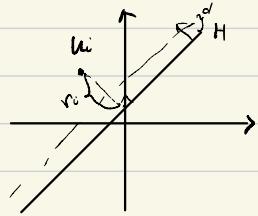


$$d = d_S(H, U_0) = \frac{|g(\underline{u})|}{\|\underline{w}\|} = \frac{b}{\|\underline{w}\|}$$

$$g(\underline{u}) = \underline{w}^T \underline{u} + w_0$$

Augmented reflected

$g(\underline{u}) = \underline{w}^T \underline{u}_i + w_0 > 0 \Rightarrow \underline{u}_i$ is correctly classified.
 $\therefore \sum_i (\underline{w}^T \underline{u}_i + w_0) > 0 \quad \forall i \Rightarrow$ all data points correctly classified.



$$\text{Distance: } r_i = d_s(H, \underline{u}_i) = \frac{\|\underline{w}\|}{\|\underline{w}\|} > 0$$

want: $r_i \geq d > 0$

$\therefore \text{Require (1)}$

Maximum margin classifier

$$\left\{ \begin{array}{l} \frac{\sum_i (\underline{w}^T \underline{u}_i + w_0)}{\|\underline{w}\|} \geq d > 0 \quad \forall i \\ \text{with maximum possible } d \end{array} \right.$$

assume L.S. in \underline{u} space
 ↵ for training dataset

II. SVM Learning

$$\text{RHS of (1): } d = \frac{b}{\|\underline{w}\|}$$

$\max d \Rightarrow \max b$ & $\min \|\underline{w}\|$

if Fix b then minimize $\|\underline{w}\|$ for given b

$\therefore (1) \Rightarrow \text{Require } \sum_i (\underline{w}^T \underline{u}_i + w_0) \geq 1 \quad \forall i, \text{ with minimal } \|\underline{w}\|$

Use Lagrange Optimization

$$(2) \left\{ \begin{array}{l} \text{Minimize } J(\underline{w}) = \|\underline{w}\|^2 \\ \text{Subject to constraints: } \sum_i (\underline{w}^T \underline{u}_i + w_0) - 1 \geq 0 \quad \forall i \end{array} \right.$$

Note: Use λ instead of μ for Lagrange multipliers

$$(2'): L(\underline{w}, w_0, \lambda) = \frac{1}{2} \|\underline{w}\|^2 - \sum_{i=1}^N \lambda_i [\sum_i (\underline{w}^T \underline{u}_i + w_0) - 1]$$

$$\lambda_i \geq 0, \quad \forall i$$

$$\lambda_i [\sum_i (\underline{w}^T \underline{u}_i + w_0) - 1] = 0 \quad \forall i$$

$$\sum_i (\underline{w}^T \underline{u}_i + w_0) - 1 \geq 0 \quad \forall i$$

Primal form of lagrangian

\underline{w}^* → solution

Primal Form:

Dual form

Easier to use if express its dual representation.

Solve: $\left\{ \begin{array}{l} \nabla_{\underline{w}} L(\underline{w}, w_0, \lambda) = 0 \\ \text{and } \nabla_{\lambda} L(\underline{w}, w_0, \lambda) = 0 \end{array} \right\}$ for \underline{w}^* \Rightarrow to get $\text{Min}(\underline{w})$ & $\text{Min}(w_0)$

(2'')

$$L_D(\lambda) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j \sum_i \sum_j \underline{u}_i^T \underline{u}_j + \sum_{i=1}^N \lambda_i$$

$$\text{with: } \sum_{i=1}^N \lambda_i z_i = 0$$

$$\lambda_i \geq 0, \lambda_i [\sum_i (\underline{w}^T \underline{u}_i + w_0) - 1] = 0 \quad \forall i$$

$$\underline{w}^* = \sum_{i=1}^N \lambda_i \underline{z}_i \underline{u}_i$$

$$\sum_i (\underline{w}^T \underline{u}_i + w_0) - 1 \geq 0$$

← DUAL FORM

$$L(\underline{\lambda}) = L(w^*, w_0, \underline{\lambda}) = \text{L optimized w.r.t. } w, w_0$$

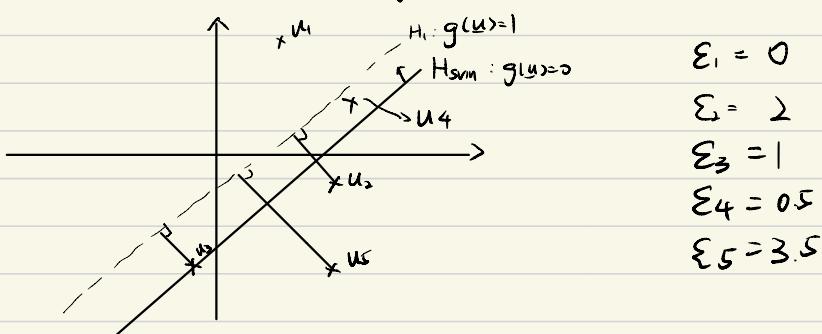
Comments:

1. $L(\underline{\lambda}) \Rightarrow$ only need to optimize w.r.t. $\underline{\lambda}$: yields solution $\underline{\lambda}^*$
2. In $L(\underline{\lambda})$, w appears in terms of what form?
 $\lambda_i z_j (z_i^T w)^2 (z_j^T w) \rightarrow$ form of a kernel
 \Rightarrow Can define nonlinear transformation by a choice of kernel $k(z_i, z_j)$

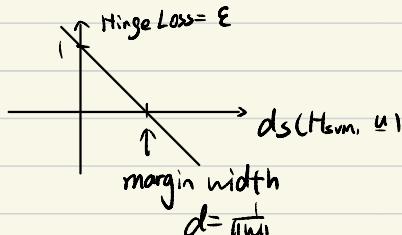
$$\begin{aligned} \nabla_w L &= \frac{1}{2} \cdot 2w - \sum_{i=1}^N \lambda_i z_i z_i^T w_i = 0 \\ \therefore w^* &= \sum_{i=1}^N \lambda_i z_i z_i^T w_i = 0 \\ \therefore \frac{\partial L}{\partial w_0} &= -\sum_{i=1}^N \lambda_i z_i w^T w_i + \lambda_i z_i w_0 + \lambda_i = \sum_{i=1}^N \lambda_i z_i = 0 \quad \text{---} \\ \text{substitute (2) into (1)} \\ L_D &= \frac{1}{2} \left\| \sum_{i=1}^N \lambda_i z_i w_i \right\|^2 - \sum_{i=1}^N \lambda_i z_i w^T w_i + \lambda_i z_i w_0 - \lambda_i \\ &= \frac{1}{2} \left\| \sum_{i=1}^N \lambda_i z_i w_i \right\|^2 - \sum_{i=1}^N \lambda_i z_i \left(\sum_{j=1}^N \lambda_j z_j w_j^T w_i \right) + \sum_{i=1}^N \lambda_i \\ &= -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j z_i z_j w_j^T w_i + \sum_{i=1}^N \lambda_i \end{aligned}$$

→ Steps from (2') to (2'')

- Introduce slack variables ξ_i , such that (ξ is x_0)
- $\xi_i = 0$ if w_i is on correct side of margin boundary
 - ξ_i = normalized distance to boundary, if on incorrect side of margin boundary
 - $\xi_i = 1$ if w_i is on decision boundary



Loss function:



Now, (2) becomes:

$$(3) \quad \begin{cases} \text{minimize} & J = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i \\ \text{constraints:} & z_i (w^T w_i + w_0) \geq 1 - \xi_i \quad \forall i \end{cases}$$

L becomes

$$\begin{aligned} L(w, w_0, \underline{\xi}, \underline{\lambda}) &= \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \lambda_i [z_i (w^T w_i + w_0) - 1 + \xi_i] - \sum_{i=1}^N \mu_i \xi_i \\ \lambda_i \geq 0 & \quad \lambda_i [z_i (w^T w_i + w_0) - 1 + \xi_i] = 0 \quad \forall i \\ \mu_i \geq 0 & \quad \mu_i \xi_i = 0 \quad \forall i \\ z_i (w^T w_i + w_0) - 1 + \xi_i \geq 0, \xi_i \geq 0 & \quad \forall i \end{aligned}$$

Primal

From L, we can derive L_D :
 $\Rightarrow L_D(\lambda)$ has same L_D for separable case, except
 constraint $\lambda_i \geq 0 \forall i \Rightarrow 0 \leq \lambda_i \leq C$

$$\therefore L_D(\lambda) = \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j \underbrace{\sum_{k,l} z_k z_l u_k^T u_l}_{z_k z_l \phi(x_k) \phi(x_l)}$$

$$\therefore k(z_k z_l) = \phi(z_k) \phi(z_l)$$

\Rightarrow we can

- ① choose a kernel function
- ② substitute kernel function

$$\begin{aligned} \nabla_w L &= \frac{1}{2} \cdot 2w - \sum_{i=1}^N \lambda_i z_i u_i = 0 \\ \therefore w^* &= \sum_{i=1}^N \lambda_i z_i u_i \quad \text{①} \\ \therefore \frac{\partial L}{\partial \lambda_i} &= -\sum_{j=1}^N 2z_j w^T u_j + \lambda_i 2z_i w_0 + \lambda_i = \sum_{j=1}^N z_j z_i = 0 \quad \text{②} \\ \frac{\partial L}{\partial \lambda_j} &= C \sum_{i=1}^N z_i \lambda_i + C - \sum_{i=1}^N z_i \lambda_i - \lambda_j - \lambda_i = -w \\ &= C - \lambda_j - \lambda_i = 0 \quad \text{③} \\ \text{Substitute } \text{②} \text{ & } \text{③} \\ L &= -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j z_i z_j u_i^T u_j + \frac{1}{2} \lambda_i^2 + C \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \lambda_i z_i \\ &\stackrel{\text{①}}{=} -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j z_i z_j u_i^T u_j + \frac{1}{2} \lambda_i^2 + C \lambda_i - \lambda_i z_i \\ \text{with } \begin{cases} \text{① } C \lambda_i - \lambda_i = 0 \\ \text{② } \lambda_i \geq 0 \\ \text{③ } 0 \leq \lambda_i \leq C \\ \text{④ } \lambda_i \in \mathbb{R} \\ \text{⑤ } \lambda_i [z_i(u_i^T w + w_0) - \epsilon] = 0 \end{cases} &= 0 \end{aligned}$$

SVR: Support Vector Regression

SVM adapted to regression problems $\begin{cases} x \rightarrow \text{NL mapping} \\ \phi(x) \text{ or } k(x, x') \rightarrow u \rightarrow \boxed{\text{linear regressor}} \end{cases} \rightarrow \hat{y}$

Motivation: to have a regression model suitable to sparse data
 e.g. high D' low N

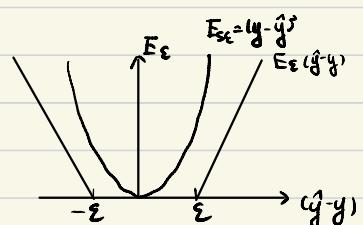
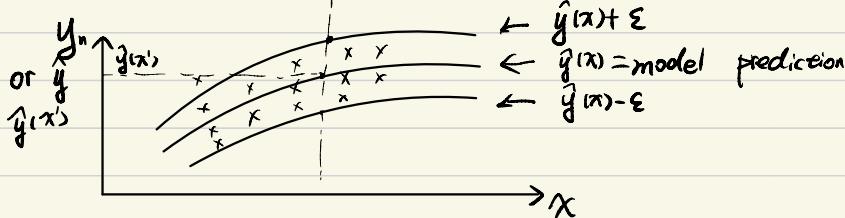
In Ridge Regression, we used a criterion fcn: $\text{MSE} + \ell^2$ regularizer

$$\text{JRR}(w) = \frac{1}{n} \sum_{i=1}^n [\hat{y}(x_i) - y_i]^2 + \lambda \|w\|_2^2, \lambda \geq 0$$

For SVR:
 use " ϵ -insensitive" loss instead of squared-error loss

ϵ -insensitive-loss E_ϵ :

$$E_\epsilon(\hat{y} - y) = \begin{cases} 0 & \text{if } |\hat{y} - y| < \epsilon \\ |\hat{y} - y| - \epsilon & \text{if } |\hat{y} - y| \geq \epsilon \end{cases}$$



- 0 error for predictions $\hat{y}(x_i)$ within ϵ of known output y_i , means learning algorithm is less likely to overfit to noise or small variations in data.