

Machine Learning I: Supervised Methods

B. Keith Jenkins

Announcements



- Slido event code: **1761012**
- Homework 5 is due Friday.
- Midterm exam is 3/6/2024, 4:00 - 5:50 PM.
 - On paper
 - Closed book and no internet
 - One formula sheet allowed (2 sided)
 - Basic scientific calculator allowed (stand-alone HW based device)

Reading

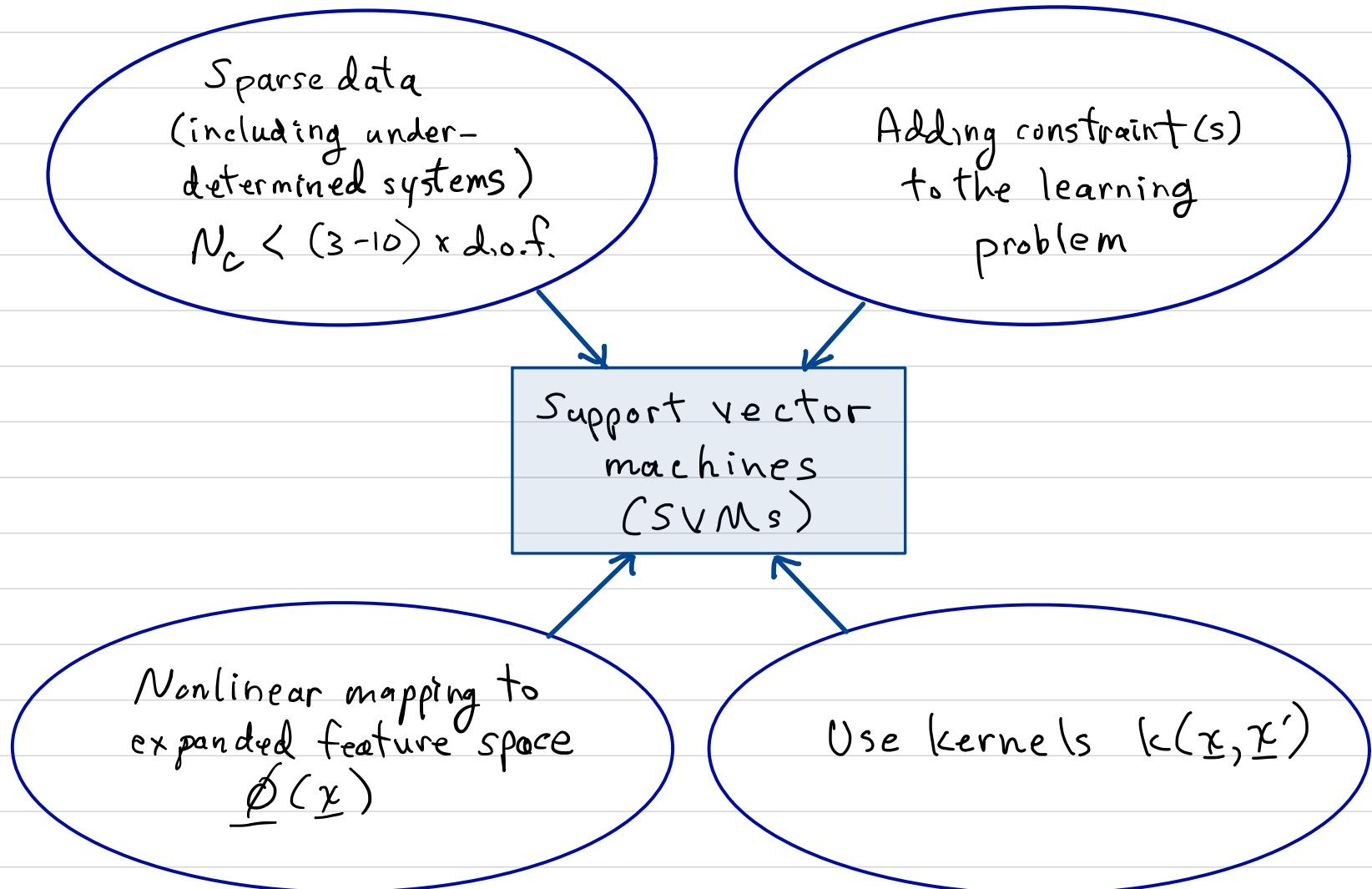
- Bishop 7.1.4 - Support Vector Regression

Today's lecture

- Support Vector Machines
 - Maximum margin classifier
 - Primal and dual forms
 - Linearly separable data
 - Not linearly separable data
 - Slack variables and hinge loss
- Support vector regression (part 1)
 - Epsilon-insensitive loss

SUPPORT VECTOR MACHINES (SVMs)

[Bishop 7.1]



Concept

→ $C=2$ classes.

→ Non-augmented notation, space

$$(i) \text{ Nonlinear mapping: } \underline{x} \xrightarrow{\text{N.L. mapping}} \underline{u} = \phi(\underline{x})$$

→ Assume data is linearly separable in \underline{u} -space (nonseparable case - later)

(ii) Optimized version of using a margin

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

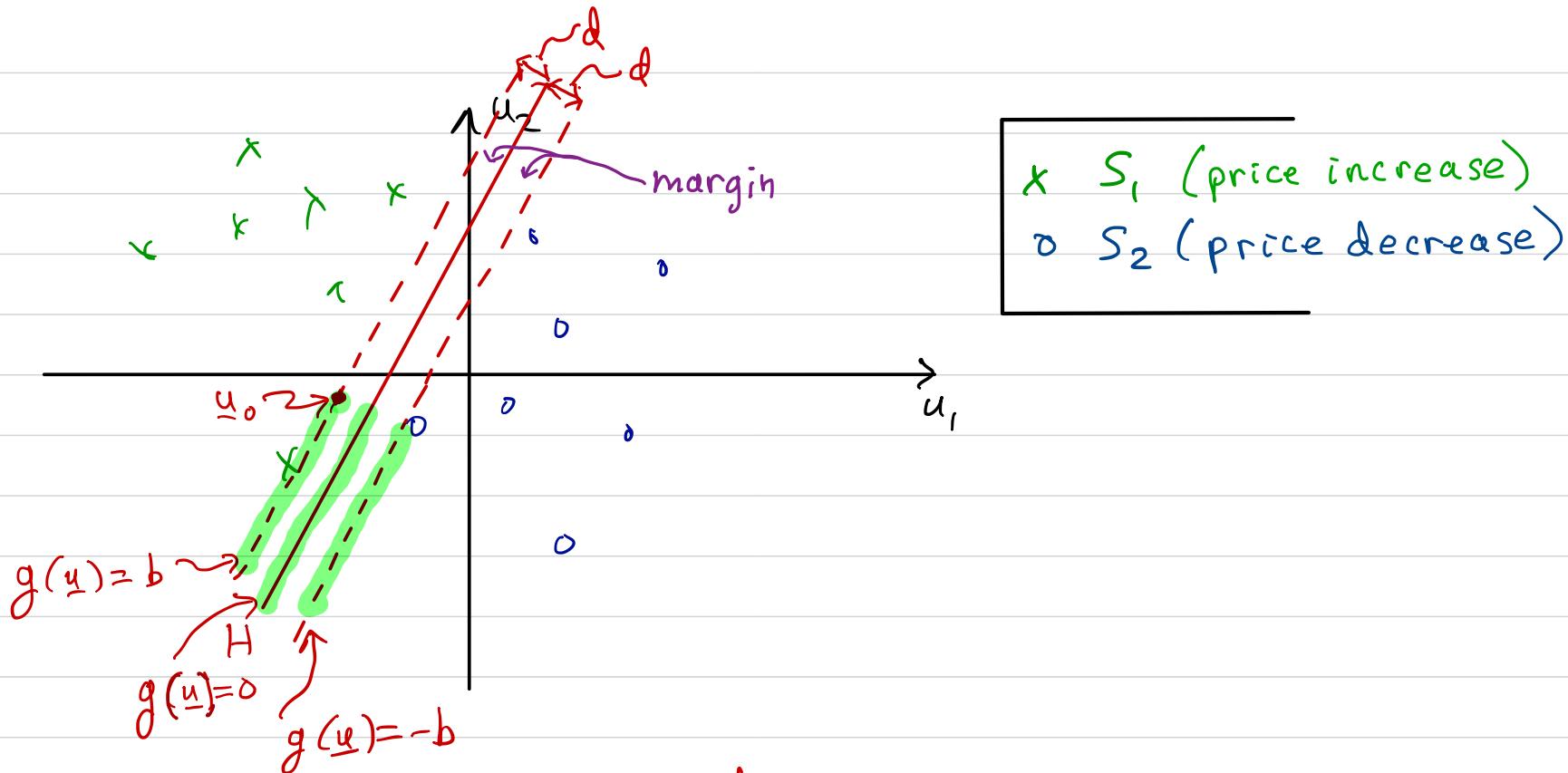
$$g'(\underline{u}) = \underline{w}^T \underline{u} + w_0' \geq b' \geq 0 \quad \text{for } \underline{u} \in S_1$$

$\phi(\underline{x})$ or $k(\underline{x}, \underline{x}')$

→ We will drop primes from hereon, and work in \underline{u} space.

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

↗ for correct classification of all \underline{u}_n .

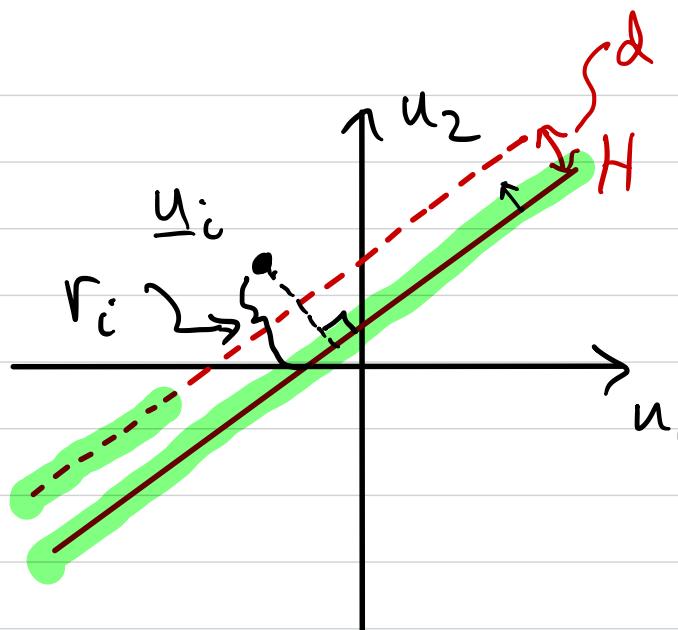


$$d = d_s(H, \underline{u}_0) = \frac{g(\underline{u}_0)}{\|\underline{w}\|} = \frac{b}{\|\underline{w}\|}$$

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

(S) augm. refl. $\Rightarrow g(\underline{u}) = \underline{w}^{(+)}{}^T z_i \underline{u}_i^{(+)} > 0 \Rightarrow \underline{u}_i$ is correctly classified.

[
 $\therefore z_i (\underline{w}^T \underline{u}_i + w_0) > 0 \quad \forall i \Rightarrow$ all data pts. correctly classified.
 —]



Distance:

$$d_s(H, \underline{u}_i) = r_i = \frac{z_i(\underline{w}^T \underline{u}_i + w_0)}{\|\underline{w}\|} > 0$$

$$\Rightarrow \text{Want: } r_i \geq d > 0$$

\therefore Require:

$$(1) \left\{ \begin{array}{l} \frac{z_i(\underline{w}^T \underline{u}_i + w_0)}{\|\underline{w}\|} \geq d > 0 \\ \text{with maximum possible } d. \end{array} \right. \quad \forall i$$

SVM Learning

(Still assuming data is linearly separable in w space.)

Note: $\underline{w} \rightarrow \alpha \underline{w}$, $w_0 \rightarrow \alpha w_0$, ANY $\alpha > 0$

\Rightarrow LHS of (1) doesn't change.

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

$\max. d \Rightarrow \max. b$, or $\min. \|\underline{w}\|$.

$\Rightarrow \min. \|\underline{w}\|$ for a given b .

and fix b by letting $b=1$.

⑤ $\therefore (1) \Rightarrow$

$$\text{Require } z_i (\underline{w}^T \underline{y}_i + w_0) \geq 1 \quad \forall i$$

with minimal $\|\underline{w}\|$.

Use Lagrange Optimization

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

Note: use λ instead of μ for Lagrange multipliers

$$\boxed{\begin{aligned} L(\underline{w}, w_0, \lambda) &= \frac{1}{2} \|\underline{w}\|^2 - \sum_{i=1}^N \lambda_i [z_i (\underline{w}^T \underline{u}_i + w_0) - 1] \\ (2') \quad \begin{cases} \lambda_i \geq 0 \quad \forall i \\ \lambda_i [z_i (\underline{w}^* T \underline{u}_i + w_0^*) - 1] = 0 \quad \forall i \\ z_i (\underline{w}^T \underline{u}_i + w_0) - 1 \geq 0 \quad \forall i \end{cases} & \left. \right\} \text{KKT conditions} \end{aligned}}$$

↗ Primal form of Lagrangian (linearly separable case)

SVM - DUAL LAGRANGIAN FOR SEPARABLE CASE

Easier to use if express in its dual representation:

Solve: $\left\{ \begin{array}{l} \nabla_{\underline{w}} L(\underline{w}, w_0, \underline{\lambda}) = \underline{0} \\ \text{and } \frac{\partial}{\partial w_0} L(\underline{w}, w_0, \underline{\lambda}) = 0 \end{array} \right\}$ for \underline{w}^* .

Substitute into (2') \Rightarrow

(2'')

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

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

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

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

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

←
Dual form
of Lagrangian
(linearly
separable case)

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

Comments

1. $L_0(\underline{\lambda}) \Rightarrow$ only need to optimize w.r.t. $\underline{\lambda}$; yields solution $\underline{\lambda}^*$.
2. In $L_0(\underline{\lambda})$, \underline{u} appears in terms of what form?

$$\lambda_i \lambda_j (z_i \underline{u}_i)^T (z_j \underline{u}_j) \rightarrow \text{form of a kernel!}$$

\Rightarrow can define nonlinear transformation by a choice of kernel $k(\underline{x}_i, \underline{x}_j)$.

Implementation

For simple problems, can be solved algebraically.

For more complex problems, a variety of numerical techniques can be used, including:

1. Quadratic programming (Optimizes quadratic fn. with constraints)

2. Projected gradient descent / sub-gradient methods
(e.g., Pegasos).

SVM Learning without Linearly Separable Assumption

→ Introduce slack variables ξ_i , such that: (ξ is x_i)

$\xi_i = 0$ if y_i is on correct side of margin boundary.

$\xi_i = \text{normalized distance to margin boundary, if}$
 $\text{on incorrect side of margin boundary.}$

$\xi_i = 1$ if y_i is on decision boundary.

F

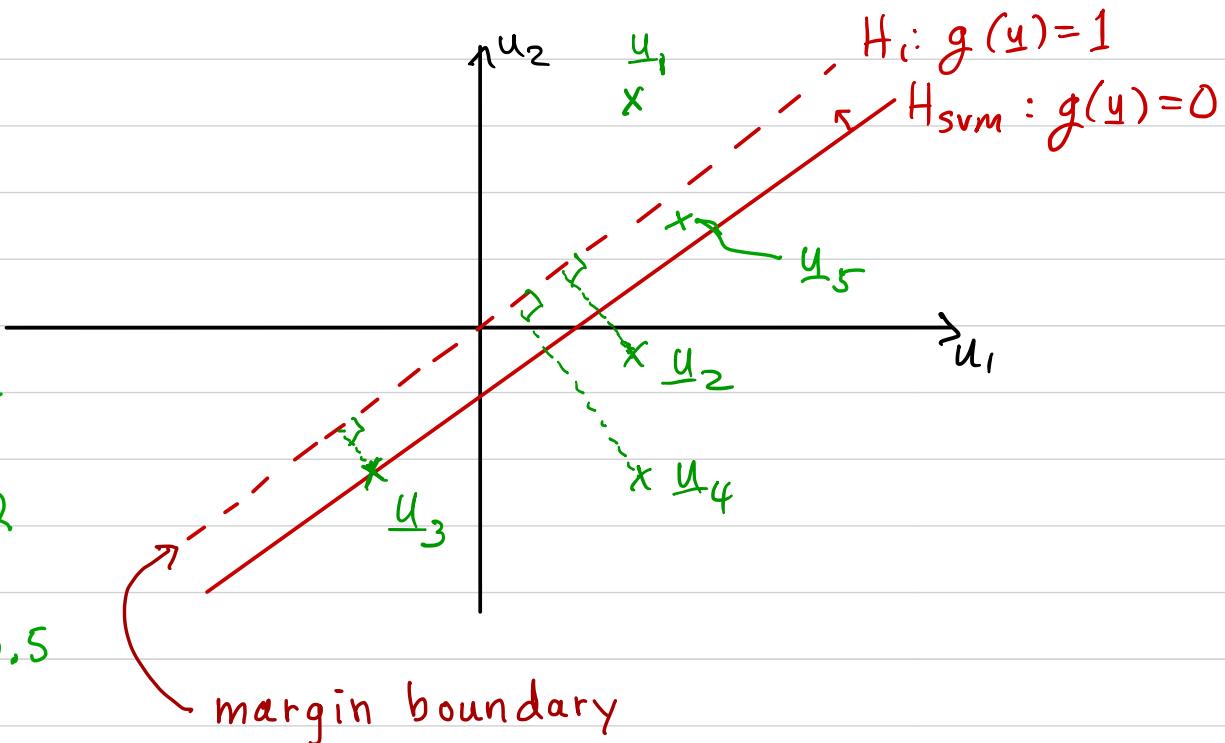
$$\xi_1 = 0$$

$$\xi_3 = 1$$

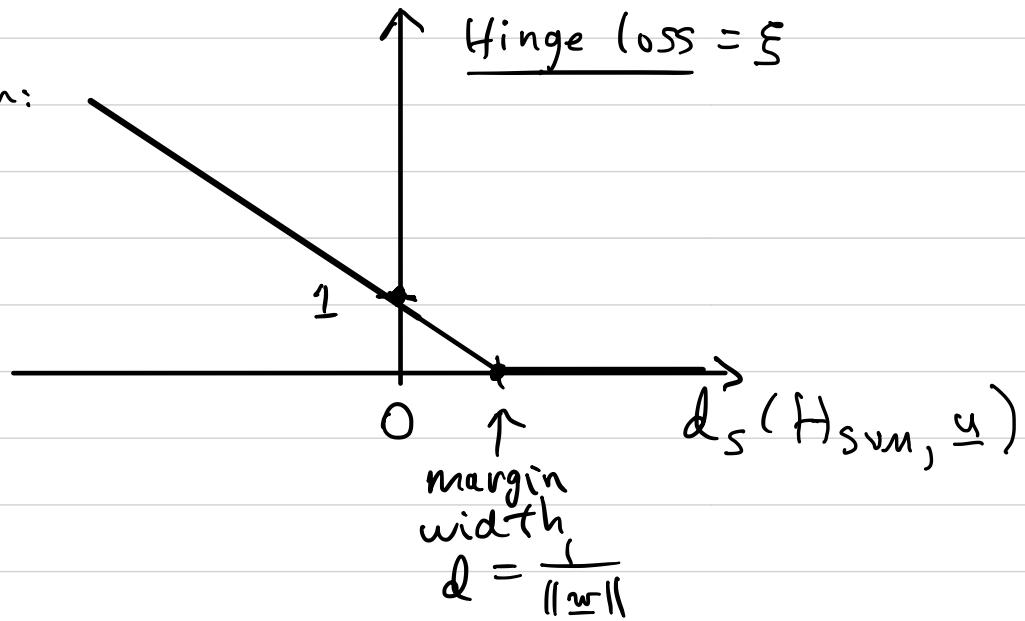
$$\xi_2 = 2$$

$$\xi_5 = 0.5$$

$$\xi_4 \approx 3.5$$



\Rightarrow Loss function:



Now, (2) becomes:

$$(3) \quad \left\{ \begin{array}{l} \text{Minimize } J = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i \\ \text{new term} \\ \text{slack-variable parameter.} \end{array} \right.$$

Constraints:

$$z_i(w^\top y_i + w_0) \geq 1 - \xi_i \quad \forall i$$

new term

Generates slightly different L (w, w_0, λ) and $L_D(\lambda)$, and constraints.

SVM With Slack Variables

L BECOMES:

(3')

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

} KKT
cond-
itions

Primal L (no linearly sep. assumption)

From L , we can derive L_D :

$\Rightarrow L_D(\underline{\lambda})$ has the same form as L_D for the separable case, except that:
 constraint $\lambda_i \geq 0 \quad \forall i$
 becomes $0 \leq \lambda_i \leq C \quad \forall i$.

BACK TO NONLINEAR MAPPING

$$L_D(\lambda) = \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j z_i z_j u_i^T u_j$$

$z_i z_j \underline{\phi}^T(x_i) \underline{\phi}(x_j)$

| $k(z_i, x_j) = \underline{\phi}^T(x_i) \underline{\phi}(x_j)$ is the kernel function

$$L_D(\lambda) = \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j z_i z_j k(x_i, x_j)$$

\Rightarrow We can:

1. Pick a kernel function
2. Use kernel substitution

to implement the nonlinear transformation,
or to use a similarity function appropriate for the
application.

Support Vector Regression (SVR) [Bishop 7.1.4] (non-augmented notation)

SVM adapted to regression problems.

Motivation: to have a regression model suitable to sparse data
(e.g., high D , relatively low N).

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

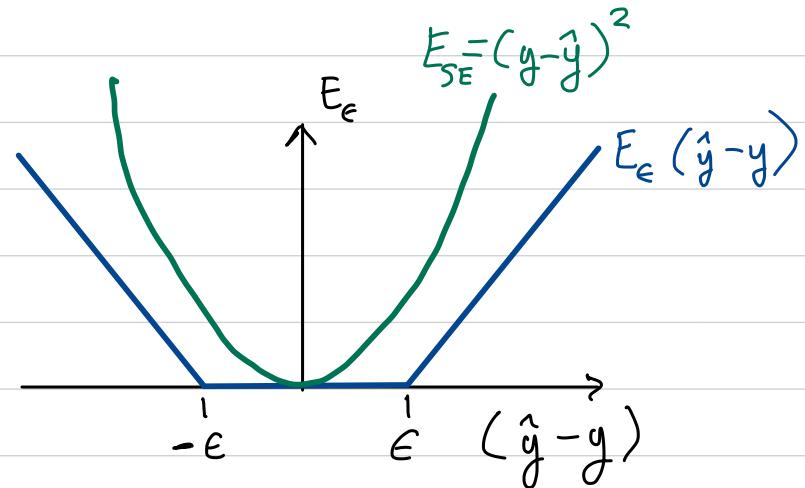
$$J_{RR}(\underline{w}) = \underbrace{\frac{1}{N} \sum_{n=1}^N [\hat{y}(x_n) - y_n]^2}_{\text{use } \epsilon\text{-insensitive loss}} + \underbrace{\lambda \|\underline{w}\|_2^2}_{\text{keep } \lambda^2 \text{ regularizer}}, \lambda \geq 0$$

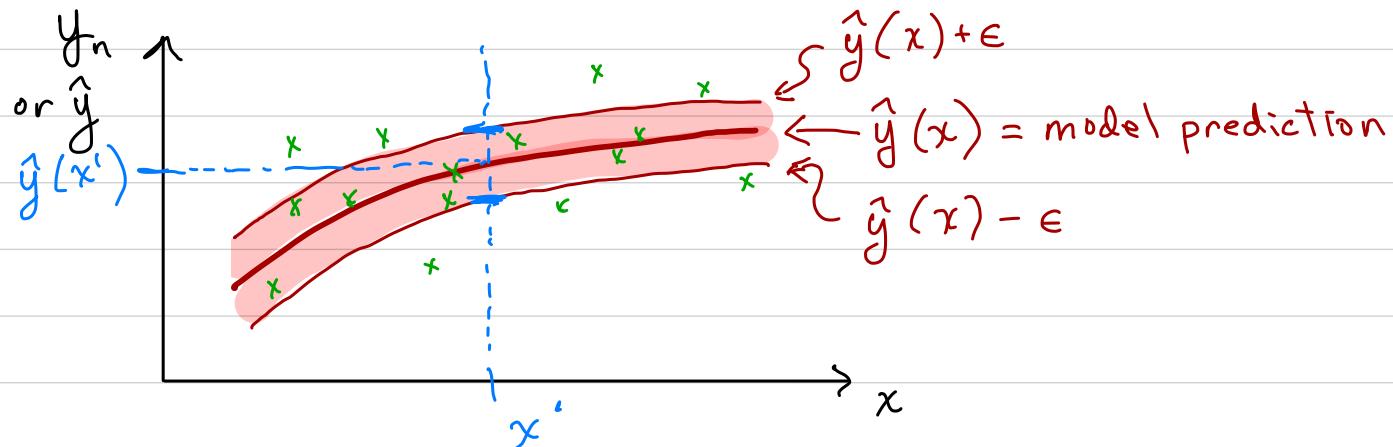
For SVR:

use " ϵ -insensitive" loss
instead of squared-error loss

ϵ -insensitive loss E_ϵ :

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





Any data points
in the shaded region
have $E_\epsilon = 0$.

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