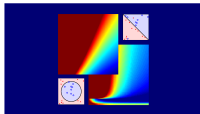


# Machine Learning Techniques (機器學習技巧)



## Lecture 6: Kernel Models for Regression

Hsuan-Tien Lin (林軒田)

htlin@csie.ntu.edu.tw

Department of Computer Science  
& Information Engineering

National Taiwan University  
(國立台灣大學資訊工程系)



# Agenda

## Lecture 6: Kernel Models for Regression

- Kernel Ridge Regression
- Support Vector Regression Primal
- Support Vector Regression Dual
- Summary of Kernel Models

## Recall: Representer Theorem

for any L2-regularized linear model

$$\min_{\mathbf{w}} \quad \frac{\lambda}{N} \mathbf{w}^T \mathbf{w} + \frac{1}{N} \sum_{n=1}^N \text{err}(y_n, \mathbf{w}^T \mathbf{z}_n)$$

optimal  $\mathbf{w}_* = \sum_{n=1}^N \beta_n \mathbf{z}_n$ .

—any L2-regularized linear model can be **kernelized!**

regression with squared error

$$\text{err}(y, \mathbf{w}^T \mathbf{z}) = (y - \mathbf{w}^T \mathbf{z})^2$$

—analytic solution for linear/ridge regression

**analytic solution** for **kernel ridge regression?**

# Kernel Ridge Regression Problem

solving ridge regression  $\min_{\mathbf{w}} \frac{\lambda}{N} \mathbf{w}^T \mathbf{w} + \frac{1}{N} \sum_{n=1}^N (y_n - \mathbf{w}^T \mathbf{z}_n)^2$

yields optimal solution  $\mathbf{w}_* = \sum_{n=1}^N \beta_n \mathbf{z}_n$

with out loss of generality, can solve for optimal  $\beta$  instead of  $\mathbf{w}$

$$\min_{\beta} \frac{\lambda}{N} \underbrace{\sum_{n=1}^N \sum_{m=1}^N \beta_n \beta_m K(\mathbf{x}_n, \mathbf{x}_m)}_{\text{regularization of } \beta \text{ on } K\text{-based regularizer}} + \frac{1}{N} \underbrace{\sum_{n=1}^N \left( y_n - \sum_{m=1}^N \beta_m K(\mathbf{x}_n, \mathbf{x}_m) \right)^2}_{\text{linear regression of } \beta \text{ on } K\text{-based features}}$$

$$= \frac{\lambda}{N} \beta^T \mathbf{K} \beta + \frac{1}{N} \left( \beta^T \mathbf{K}^T \mathbf{K} \beta - 2 \beta^T \mathbf{K}^T \mathbf{y} + \mathbf{y}^T \mathbf{y} \right)$$

kernel ridge regression:

use **representer theorem** for kernel trick on **ridge regression**

## Solving Kernel Ridge Regression

$$E_{\text{aug}}(\boldsymbol{\beta}) = \frac{\lambda}{N} \boldsymbol{\beta}^T \mathbf{K} \boldsymbol{\beta} + \frac{1}{N} \left( \boldsymbol{\beta}^T \mathbf{K}^T \mathbf{K} \boldsymbol{\beta} - 2 \boldsymbol{\beta}^T \mathbf{K}^T \mathbf{y} + \mathbf{y}^T \mathbf{y} \right)$$

$$\nabla E_{\text{aug}}(\boldsymbol{\beta}) = \frac{2}{N} \left( \lambda \mathbf{K}^T \mathbf{I} \boldsymbol{\beta} + \mathbf{K}^T \mathbf{K} \boldsymbol{\beta} - \mathbf{K}^T \mathbf{y} \right) = \frac{2}{N} \mathbf{K}^T \left( (\lambda \mathbf{I} + \mathbf{K}) \boldsymbol{\beta} - \mathbf{y} \right)$$

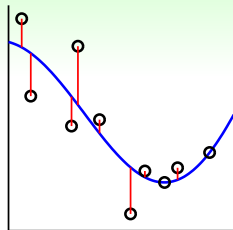
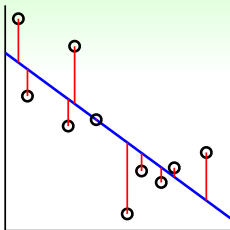
want  $\nabla E_{\text{aug}}(\boldsymbol{\beta}) = \mathbf{0}$ : one analytic solution

$$\boldsymbol{\beta} = (\lambda \mathbf{I} + \mathbf{K})^{-1} \mathbf{y}$$

- $(\cdot)^{-1}$  always exists for  $\lambda > 0$ , because  $\mathbf{K}$  positive semi-definite (**Mercer's condition, remember? :-)**)
- time complexity:  $O(N^3)$  with simple **dense** matrix inversion

can now do **non-linear regression** 'easily'

# Linear versus Kernel Ridge Regression



## linear ridge regression

$$\mathbf{w} = (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- more restricted
- $O(d^3 + d^2 N)$  training;  
 $O(d)$  prediction  
— **efficient when  $N \gg d$**

## kernel ridge regression

$$\beta = (\lambda \mathbf{I} + \mathbf{K})^{-1} \mathbf{y}$$

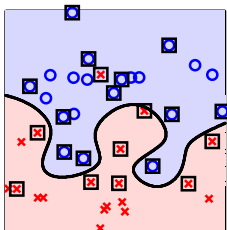
- **more flexible** with  $K$
- $O(N^3)$  training;  
 $O(N)$  prediction  
— hard for big data

**linear** versus **kernel**:  
trade-off between **efficiency** and **flexibility**

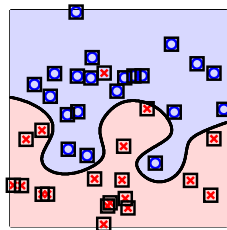
# Fun Time

# Soft-Margin SVM versus Least-Squares SVM

least-squares SVM (LSSVM)  
= **kernel ridge regression** for classification



soft-margin Gaussian SVM



Gaussian LSSVM

- LSSVM: similar boundary, **many more SVs**  
⇒ slower prediction, **dense  $\beta$  (BIG  $g$ )**
- dense  $\beta$ : LSSVM, kernel LogReg; **sparse  $\alpha$ : standard SVM**

want: **sparse  $\beta$**  like standard SVM



# Tube Regression

will consider **tube regression**

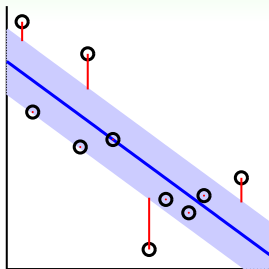
- within a **tube**: **no error**
- outside a tube: **error** by distance to tube

error measure:

$$\text{err}(y, s) = \max(0, |s - y| - \epsilon)$$

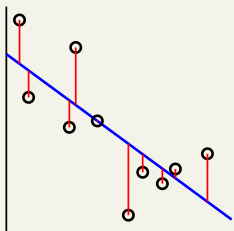
- $|s - y| \leq \epsilon$ : 0
- $|s - y| > \epsilon$ :  $|s - y| - \epsilon$

—usually called  $\epsilon$ -insensitive error with  $\epsilon > 0$

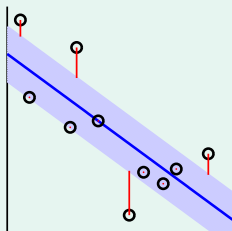


todo: **L2-regularized tube regression**  
to get **sparse  $\beta$**

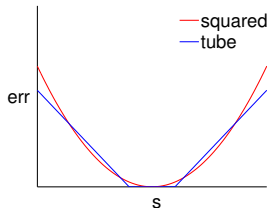
# Tube versus Squared Regression



**tube:**  $\text{err}(y, s) = \max(0, |s - y| - \epsilon)$



**squared:**  $\text{err}(y, s) = (s - y)^2$



**tube**  $\approx$  **squared** when  $|s - y|$  small  
& **less affected by outliers**

# L2-Regularized Tube Regression

$$\min_{\mathbf{w}} \quad \frac{\lambda}{N} \mathbf{w}^T \mathbf{w} + \frac{1}{N} \sum_{n=1}^N \max \left( 0, |\mathbf{w}^T \mathbf{z}_n - y| - \epsilon \right)$$

## Regularized Tube Regr.

$$\min \frac{\lambda}{N} \mathbf{w}^T \mathbf{w} + \frac{1}{N} \sum \text{tube violation}$$

- unconstrained, but **max not differentiable**
- 'representer' to kernelize, but **no obvious sparsity**

## standard SVM

$$\min \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum \text{margin vio.}$$

- not differentiable, but **QP**
- dual to kernelize, KKT conditions  $\Rightarrow$  **sparsity**

will mimic **standard SVM** derivation:

$$\min_{b, \mathbf{w}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \max \left( 0, |\mathbf{w}^T \mathbf{z}_n + b - y_n| - \epsilon \right)$$

# Standard Support Vector Regression Primal

$$\min_{b, \mathbf{w}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \max(0, |\mathbf{w}^T \mathbf{z}_n + b - y_n| - \epsilon)$$

mimicking standard SVM

$$\begin{aligned} \min_{b, \mathbf{w}, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \xi_n \\ \text{s.t.} \quad & |\mathbf{w}^T \mathbf{z}_n + b - y_n| \leq \epsilon + \xi_n \\ & \xi_n \geq 0 \end{aligned}$$

making constraints linear

$$\begin{aligned} \min_{b, \mathbf{w}, \xi^V, \xi^A} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N (\xi_n^V + \xi_n^A) \\ \text{s.t.} \quad & -\epsilon - \xi_n^V \leq \mathbf{w}^T \mathbf{z}_n + b - y_n \leq \epsilon + \xi_n^A \\ & \xi_n^V \geq 0, \xi_n^A \geq 0 \end{aligned}$$

Support Vector Regression (SVR) primal:  
 minimize regularizer +  
 upper tube violations  $\xi_n^A$  & lower violations  $\xi_n^V$

# Fun Time

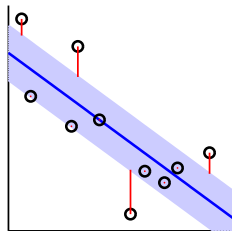
# Quadratic Programming for SVR

$$\min_{b, \mathbf{w}, \xi_n^V, \xi_n^A} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N (\xi_n^V + \xi_n^A)$$

$$\text{s.t.} \quad -\epsilon - \xi_n^V \leq \mathbf{w}^T \mathbf{z}_n + b - y_n \leq \epsilon + \xi_n^A$$

$$\xi_n^V \geq 0, \xi_n^A \geq 0$$

- parameter  $C$ : trade-off of regularization & tube violation
- parameter  $\epsilon$ : vertical tube width  
—one more parameter to choose!
- QP of  $\tilde{d} + 1 + 2N$  variables,  $2N + 2N$  constraints



next: remove dependence on  $\tilde{d}$  by  
SVR primal  $\Rightarrow$  dual?

Lagrange Multipliers  $\alpha^\wedge$  &  $\alpha^\vee$ 

objective function  $\frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N (\xi_n^\vee + \xi_n^\wedge)$

Lagrange multiplier  $\alpha_n^\vee$  for  $-\epsilon - \xi_n^\vee \leq \mathbf{w}^T \mathbf{z}_n + b - y_n$

Lagrange multiplier  $\alpha_n^\wedge$  for  $\mathbf{w}^T \mathbf{z}_n + b - y_n \leq \epsilon + \xi_n^\wedge$

## Some of the KKT Conditions

- $\frac{\partial \mathcal{L}}{\partial \mathbf{w}_i} = 0$ :  $\mathbf{w} = \sum_{n=1}^N \underbrace{(\alpha_n^\vee - \alpha_n^\wedge)}_{\beta_n} \mathbf{z}_n$

- complementary slackness:  $\alpha_n^\vee (-\epsilon - \xi_n^\vee - \mathbf{w}^T \mathbf{z}_n - b + y_n) = 0$   
 $\alpha_n^\wedge (-\epsilon - \xi_n^\wedge + \mathbf{w}^T \mathbf{z}_n + b - y_n) = 0$

standard dual can be derived  
using the same steps as Lecture 20

## SVM Dual and SVR Dual

$$\begin{aligned} \min \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \xi_n \\ \text{s.t.} \quad & y_n (\mathbf{w}^T \mathbf{z}_n + b) \geq 1 - \xi_n \\ & \xi_n \geq 0 \end{aligned}$$

$$\begin{aligned} \min \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N (\xi_n^{\wedge} + \xi_n^{\vee}) \\ \text{s.t.} \quad & 1 (\mathbf{w}^T \mathbf{z}_n + b - y_n) \leq \epsilon + \xi_n^{\wedge} \\ & 1 (y_n - \mathbf{w}^T \mathbf{z}_n + b) \leq \epsilon + \xi_n^{\vee} \\ & \xi_n^{\wedge} \geq 0, \xi_n^{\vee} \geq 0 \end{aligned}$$

$$\begin{aligned} \min \quad & \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m K(\mathbf{x}_n, \mathbf{x}_m) \\ & - \sum_{n=1}^N 1 \cdot \alpha_n \\ \text{s.t.} \quad & \sum_{n=1}^N y_n \alpha_n = 0 \\ & 0 \leq \alpha_n \leq C \end{aligned}$$

$$\begin{aligned} \min \quad & \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (\alpha_n^{\vee} - \alpha_n^{\wedge}) (\alpha_m^{\vee} - \alpha_m^{\wedge}) k_{n,m} \\ & - \sum_{n=1}^N ((\epsilon + y_n) \cdot \alpha_n^{\wedge} + (\epsilon - y_n) \cdot \alpha_n^{\vee}) \\ \text{s.t.} \quad & \sum_{n=1}^N 1 \cdot (\alpha_n^{\vee} - \alpha_n^{\wedge}) = 0 \\ & 0 \leq \alpha_n^{\wedge} \leq C, 0 \leq \alpha_n^{\vee} \leq C \end{aligned}$$

similar QP, **solvable by similar solver**



# Sparsity of SVR Solution

- $\mathbf{w} = \sum_{n=1}^N \underbrace{(\alpha_n^{\vee} - \alpha_n^{\wedge})}_{\beta_n} \mathbf{z}_n$

- complementary slackness:

$$\alpha_n^{\vee} (-\epsilon - \xi_n^{\vee} - \mathbf{w}^T \mathbf{z}_n - b + y_n) = 0$$

$$\alpha_n^{\wedge} (-\epsilon - \xi_n^{\wedge} + \mathbf{w}^T \mathbf{z}_n + b - y_n) = 0$$

- strictly within tube  $|\mathbf{w}^T \mathbf{z}_n + b - y_n| < \epsilon$   
 $\implies \alpha_n^{\wedge} = 0$  and  $\alpha_n^{\vee} = 0$   
 $\implies \beta_n = 0$
- SVs ( $\beta_n \neq 0$ ): **on or outside tube**

SVR: allows **sparse**  $\beta$

# Fun Time

# Map of Linear Models

PLA/pocket

minimize  
 $\text{err}_{0/1}$  specially

linear SVR

minimize regularized  
 $\text{err}_{\text{TUBE}}$  by QP

linear soft-margin  
SVM

minimize regularized  
 $\widehat{\text{err}}_{\text{SVM}}$  by QP

linear ridge  
regression

minimize regularized  
 $\text{err}_{\text{SQR}}$  analytically

regularized logistic  
regression

minimize regularized  
 $\text{err}_{\text{CE}}$  by GD/SGD

second row: popular in **LIBLINEAR**

# Map of Linear/Kernel Models

PLA/pocket

linear SVR

linear soft-margin  
SVM

linear ridge  
regression

regularized logistic  
regression

kernel ridge  
regression

kernel logistic  
regression

kernelized linear ridge  
regression

kernelized regularized  
logistic regression

SVM

minimize SVM dual by  
QP

SVR

minimize SVR dual by  
QP

probabilistic SVM

run SVM-transformed  
logistic regression

fourth row: popular in **LIBSVM**

# Map of Linear/Kernel Models

PLA/pocket

linear SVR

linear soft-margin  
SVM

linear ridge  
regression

regularized logistic  
regression

kernel ridge  
regression

kernel logistic  
regression

SVM

SVR

probabilistic SVM

first row: less used due to **worse performance**  
third row: less used due to **dense  $\beta$**

# Kernel Models

possible kernels:

polynomial, Gaussian, . . . , your own from Mercer's condition,

coupled with

kernel ridge  
regression

kernel logistic  
regression

SVM

SVR

probabilistic SVM

powerful extension of linear models  
— *with great power comes great responsibility*  
in **Spiderman, remember? :-)**

# Fun Time

# Summary

## Lecture 6: Kernel Models for Regression

- Kernel Ridge Regression  
**representer theorem on RidgeReg**
- Support Vector Regression Primal  
**minimize regularized tube errors**
- Support Vector Regression Dual  
**a QP similar to SVM**
- Summary of Kernel Models  
**with great power comes great responsibility**