



Kernel Models

MLI

Module 1, Lecture 5

18th May 2021

[Adriano Koshiyama]

Overview

/// Back to Feature Maps

- /// Overview
- /// Back to Ridge Regression
- /// Explicit Feature Maps
- /// Kernel Ridge Regression
- /// Kernel Trick
- /// Pros and cons of the dual representation

/// Kernels

- /// Definition
- /// Kernelization
- /// Some key facts and results
- /// A bestiary of kernels

/// Support Vector Machines

- /// Main Idea and Details
- /// Linear Support Vector Machines
- /// Kernel/Nonlinear Support Vector Machines
- /// Kernel selection and fine-tuning

Back to Feature Maps

- /// Overview
- /// Back to Ridge Regression
- /// Explicit Feature Maps
- /// Kernel Ridge Regression
- /// Kernel Trick
- /// Pros and cons of the dual representation

Overview

/// We show how a linear method such as least squares may be lifted to a (potentially) higher dimensional space to provide a nonlinear regression.

Map // Selection	Explicit	Implicit
Explicit	Basis Functions with Stepwise	Regularization
Implicit	Kernelization with Stepwise	Kernel Methods with Regularization

/// **Implicit** Selection are widely adopted in Machine Learning

/// **Explicit** Mapping is not very systematic, being more problem-dependent, but some recipes are universal (interactions, z-scores, etc.)

/// **Implicit** Mapping with **Implicit** Feature selection is **Our Focus Today!**

/// **Implicit** mapping with **Explicit** feature selection is not generally used

Back to Ridge Regression

/// If we consider the Ridge Regression case, we have

$$\min_{\beta(\lambda)} \text{MSE}[\beta(\lambda)] = \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2, \text{ for a given } \lambda \geq 0$$

/// we can rewrite this expression as

$$\min_{\beta(\lambda)} \text{MSE}[\beta(\lambda)] = \mathbf{y}^T \mathbf{y} - 2\beta^T \mathbf{X}^T \mathbf{y} + \beta^T \mathbf{X}^T \mathbf{X} \beta + \lambda \beta^T \mathbf{I}_{J_{XJ}} \beta$$

/// First order condition: $\nabla_{\beta(\lambda)} \text{MSE}[\beta(\lambda)] = 0$, then

$$-2\mathbf{X}^T \mathbf{y} + \mathbf{X}^T \mathbf{X} \beta + \lambda \mathbf{I}_{J_{XJ}} \beta = \mathbf{0} \rightarrow -2\mathbf{X}^T \mathbf{y} + (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_{J_{XJ}}) \beta = \mathbf{0}$$

/// Yielding the solution, akin to Linear Regression, but with a new param

$$\beta = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_{J_{XJ}})^{-1} \mathbf{X}^T \mathbf{y}$$

/// Second order condition: $H[\text{MSE}[\beta(\lambda)]]$ is positive semi-definite

$$H[\text{MSE}[\beta(\lambda)]] = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_{J_{XJ}})^T = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_{J_{XJ}})$$

since $\mathbf{X}^T \mathbf{X}$ is positive semi-definite, $\mathbf{I}_{J_{XJ}}$ is positive definite, the hessian matrix will be positive semi-definite for a suitable choice of λ (usually $\lambda \geq 0$)

PS: remember to scale \mathbf{X} columns (like having zero mean and unit std)

Explicit Feature Maps

- Overall, the main idea is to extend the model, including new features engineered using some **nonlinear mapping** (a.k.a. basis function)

$$y_i = \alpha^T \phi(x_i) + \varepsilon_i$$

with $\phi: S^J \rightarrow \mathcal{R}^M$, usually $M > J$; hence more coefficients have to be fitted

- Since the model is linear in the parameters α^T , we can call $\phi(x_i) = \mathbf{z}_i$ and use the usual solution $\mathbf{w} = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{y}$, with $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_n]^T$

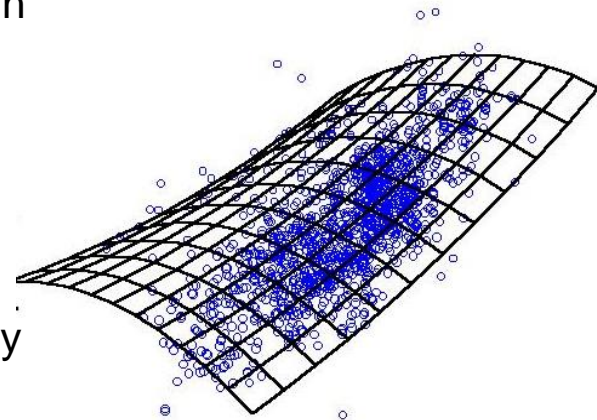
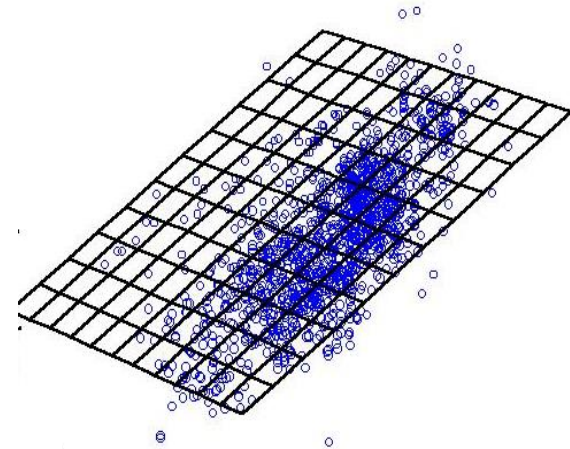
- Some commonly used feature maps, adopted in many situations

- Polynomial: $\phi: x_i \rightarrow (1, x_i, x_i^2, x_i^3, \dots)$

- Veronese/ANOVA/Interaction mapping

$\phi: x_{i1}, x_{i2} \rightarrow (1, x_{i1}, x_{i2}, x_{i1}^2, x_{i2}^2, x_{i1} x_{i2}, \dots)$

- In fact, log, sqrt, one-hot encoding, etc. and any other transformation



Kernel Ridge Regression

/// Is it possible to create polynomial features, without incurring in all the additional calculations, memory storage, etc.? **A: Yes**

/// **First step:** rewrite Ridge Regression solution by

$$X^T X \beta + \lambda I_{J_{xJ}} \beta = X^T y \Rightarrow \lambda I_{J_{xJ}} \beta = X^T y - X^T X \beta \Rightarrow \beta = \left(\frac{1}{\lambda} \right) X^T [y - X \beta]$$

hence, β can be expressed as a linear combination of the columns of X^T

$$\beta = X^T \alpha$$

with $\alpha \in \mathbb{R}^n$ being the dual coefficient vector; further manipulation show

$$\alpha = \left(\frac{1}{\lambda} \right) [y - X \beta] \Rightarrow \alpha \lambda = [y - X X^T \alpha] \Rightarrow (X X^T + \lambda I_{n \times n}) \alpha = y$$

/// therefore, $\alpha = (X X^T + \lambda I_{n \times n})^{-1} y$, implying that α

/// minimizes $MSE[\alpha(\lambda)] = \|y - X X^T \alpha\|_2^2 + \lambda \alpha^T X X^T \alpha$

/// written in terms of inner products of rows ($X X^T$) not columns ($X^T X$)

Kernel Trick

/// **Second step:** since α is written in terms of inner products of rows

$$XX^T = \begin{bmatrix} \langle x_1, x_1 \rangle & \cdots & \langle x_1, x_n \rangle \\ \vdots & \ddots & \vdots \\ \langle x_n, x_1 \rangle & \cdots & \langle x_n, x_n \rangle \end{bmatrix}$$

/// if, we consider a feature map of x_i , $\phi: S^J \rightarrow \mathcal{R}^M$, we can rewrite the Gram-Matrix XX^T after feature mapping by

$$. = \begin{bmatrix} \langle \phi(x_1), \phi(x_1) \rangle & \cdots & \langle \phi(x_1), \phi(x_n) \rangle \\ \vdots & \ddots & \vdots \\ \langle \phi(x_n), \phi(x_1) \rangle & \cdots & \langle \phi(x_n), \phi(x_n) \rangle \end{bmatrix} = \begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{bmatrix} = K$$

/// where $k(x_k, x_i) := \langle \phi(x_k), \phi(x_i) \rangle$ is a (**kernel**) function that **replicates** a certain **feature mapping plus an inner product**!

/// Example: consider that $x, z \in \mathbb{R}^2$ (think as two rows of a 2-d dataset) and $k(x, z) = \langle x, z \rangle^2$, then $k(x, z) = x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 x_2 z_1 z_2 = \cdots$

$$\dots = [x_1^2, x_2^2, \sqrt{2}x_1 x_2]^T [z_1^2, z_2^2, \sqrt{2}z_1 z_2] = \langle \phi(x), \phi(z) \rangle$$

for a map $\phi(x) = [x_1^2, x_2^2, \sqrt{2}x_1 x_2]$, close to a second degree polynomial

/// hence, $\langle \phi(x), \phi(z) \rangle$ takes roughly $O(n^2)$ whilst $k(x, z)$ takes $O(n)$

Dual Representation – Pros and Cons

/// Primal Form: Computation with Basis Functions

/// **Training:** $O(J^p)$ for feature mapping and $O(nJ^2 + J^3)$ for estimation

/// Dataset $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$

/// Basis functions: ϕ_1, \dots, ϕ_k where $\phi: \mathbb{R}^J \rightarrow \mathbb{R}$

/// Feature map: $\boldsymbol{\phi}(\mathbf{x}) = [\phi_1(\mathbf{x}), \dots, \phi_k(\mathbf{x})]$

/// Solution: $\boldsymbol{\beta} = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{y}$, with $\mathbf{Z} = [\boldsymbol{\phi}(\mathbf{x}_1), \dots, \boldsymbol{\phi}(\mathbf{x}_n)]^T$

/// **Prediction:** $\hat{y}(\mathbf{x}_i) = \boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_i)$, taking $O(k)$ operations

/// Dual Form: Implicit Feature Map

/// **Training:** $O(Jn)$ to create \mathbf{K} and $O(Jn^2 + n^3)$ for estimation

/// Dataset $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$

/// Kernel Function: $k: \mathbb{R}^J \times \mathbb{R}^J \rightarrow \mathbb{R}$

/// Kernel Matrix: $\mathbf{K} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \cdots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}$

/// Solution: $\boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I}_{n \times n})^{-1} \mathbf{y}$

/// **Prediction:** $\hat{y}(\mathbf{v}) = \sum_i \alpha_i k(\mathbf{x}_i, \mathbf{v})$, taking $O(nJ)$ operations

Kernels

- /// Definition
- /// Kernelization
- /// Some key facts and results
- /// A bestiary of kernels

Definition

/// Kernel

/// A function $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is called a kernel on \mathcal{X} if there exists a Hilbert space \mathcal{F} and a map $\phi: \mathcal{X} \rightarrow \mathcal{F}$, such that $k(x, z) = \langle \phi(x), \phi(z) \rangle_{\mathcal{F}}$

/// In our more usual terminology

/// $\phi: \mathcal{X} \rightarrow \mathcal{F}$ is called a **feature map**

/// \mathcal{F} is called a **feature space**

/// Hence, a kernel function is a composition of "feature map" with scalar product. Our previous example, $x, z \in \mathbb{R}^2$, we have:

/// $k(x, z) = \langle x, z \rangle^2$ (homogenous 2-degree polynomial kernel)

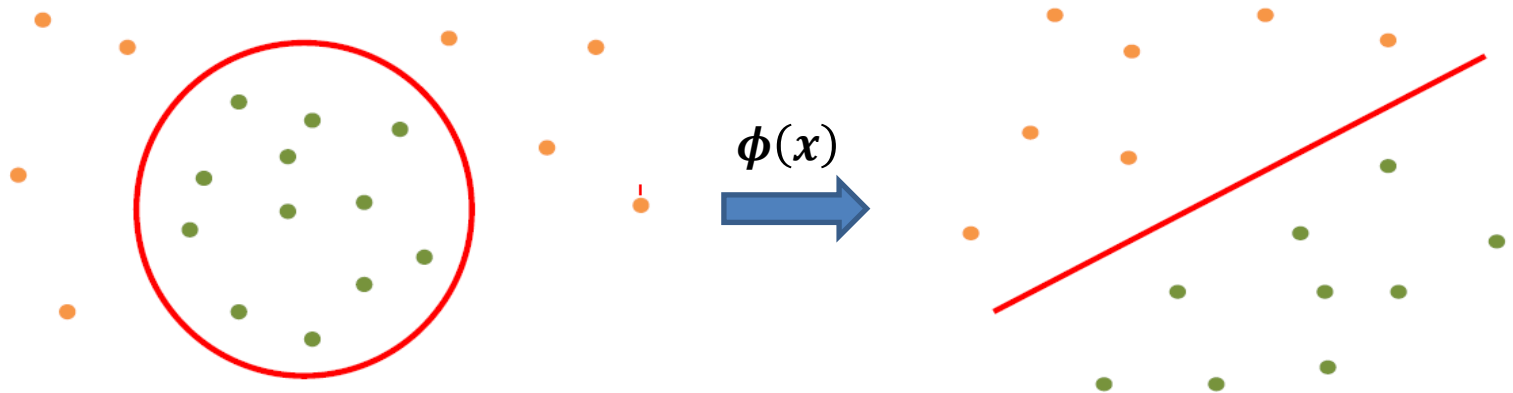
/// $\phi(x) = [x_1^2, x_2^2, \sqrt{2}x_1x_2]$ or $\phi(x) = [x_1^2, x_2^2, x_1x_2, x_1x_2]$

is one of the most common instantiations of a kernel function

/// **Note that a kernel function is not always associated with a unique feature map!**

Kernelization

- /// Kernels allow to make linear algorithms work on non-linear features, by replicating a very complicated feature mapping by a kernel function over the raw features
- /// Example: Support Vector Machine: $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n, y_i \in \{-1, 1\}$



$$f(\mathbf{x}) = \text{sgn}(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})$$

$\boldsymbol{\beta} = \sum_i \alpha_i \mathbf{x}_i y_i$ solves a QP involving $\mathbf{X}^T \mathbf{X}$, with α_i representing “dual variables”

$$\begin{aligned} f(\mathbf{z}) &= \text{sgn} \left(\beta_0 + \sum_i \alpha_i k(\mathbf{x}_i, \mathbf{z}) \right) \\ &= \boldsymbol{\alpha}^T \mathbf{k}(\mathbf{X}, \mathbf{z}) \end{aligned}$$

$\boldsymbol{\alpha}$ solves a QP involving kernel matrix \mathbf{K}

Key facts and results

/// What functions are kernels?

/// Given a function $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, which properties of k guarantee that there exists a Hilbert space \mathcal{F} and a map $\phi: \mathcal{X} \rightarrow \mathcal{F}$, such that $k(x, z) = \langle \phi(x), \phi(z) \rangle_{\mathcal{F}}$?

/// **Subtle note:** we have generalized the definition of **finite-dimensional** feature maps $\phi: \mathbb{R}^J \rightarrow \mathbb{R}^m$ to **infinite-dimensional** feature maps $\phi: \mathbb{R}^J \rightarrow \mathcal{F}$ (Hilbert space)

/// **First:** a function $h: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is **positive semidefinite** if is symmetric and the matrix H made of $H_{ixj} = h(x_i, x_j)$ for every i and j is positive semidefinite ($v^T H v \geq 0$)

/// Theorem (Positive Semidefinite Kernels)

/// A function k is positive semidefinite if and only if

$$k(x, z) = \langle \phi(x), \phi(z) \rangle, \quad x, z \in \mathbb{R}^J$$

for some feature map $\phi: \mathbb{R}^J \rightarrow \mathcal{F}$ and Hilbert space \mathcal{F}

Key facts and results

/// Theorem (Moore-Aronszajn, 1950)

/// For every positive semidefinite kernel $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, there is an unique feature map $\phi: \mathcal{X} \rightarrow \mathcal{F}$ (up to a isomorphism) into an unique Hilbert space \mathcal{F} (the so-called Reproducing Kernel Hilbert Space – RKHS)

/// **RKHS:** $\phi: x \rightarrow k(x, \cdot)$ and $\langle k(x, \cdot), f(\cdot) \rangle_{\mathcal{F}} = f(x), \forall f \in \mathcal{F}$

/// Theorem (Kimeldorf-Wahba, 1971) “Representer Theorem”

/// Let $D = \{(x_i, y_i)\}_{i=1}^n$ be data points

/// Let k be a positive semidefinite kernel with RKHS \mathcal{F}

/// Let $\mathcal{L}: \mathbb{R}^J \times \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ be any loss function

/// Let $\Omega: [0, \infty) \rightarrow \mathbb{R}$ strictly increasing penalty function

/// Then, there is a minimizer of the type $f(x) = \alpha^T k(x, \cdot)$ to

$$R(f) = \sum_i \mathcal{L}(x_i, y_i, f(x_i)) + \Omega(\|f\|_{\mathcal{F}})$$

A bestiary of kernels

/// **Linear:** $k(x, z) = x^T A z$ (usually A is a positive definite matrix)

/// **Dot-product:** $k(x, z) = f(\langle x, z \rangle)$ (usually f is a positive definite)

/// **Polynomial:** $k(x, z) = (\theta \langle x, z \rangle + c)^d$ (general formulation)

/// **Sigmoid:** $k(x, z) = \tanh(\sigma \langle x, z \rangle + c)$ (not positive definite...)

/// **Radial basis function kernels**

Gaussian/RBF Kernel

Laplacian Kernel

$$k(x, z) = \exp\left(-\frac{\|x - z\|_2^2}{2\sigma^2}\right) \quad k(x, z) = \exp\left(-\frac{\|x - z\|}{\sigma}\right)$$

/// **Proposition**

/// For positive definite kernels k_1, k_2, \dots and $c_i > 0$, any of the composition below is also a positive definite kernel

$$k_1(x, z) + c_1$$

$$c_1 k_1(x, z) + c_2 k_2(x, z)$$

$$k_1(x, z) k_2(x, z)$$

$$\frac{k_1(x, z)}{\sqrt{k_1(x, x) k_1(z, z)}}$$



$$\sum_{i=1}^{\infty} c_i \langle x, z \rangle^i$$

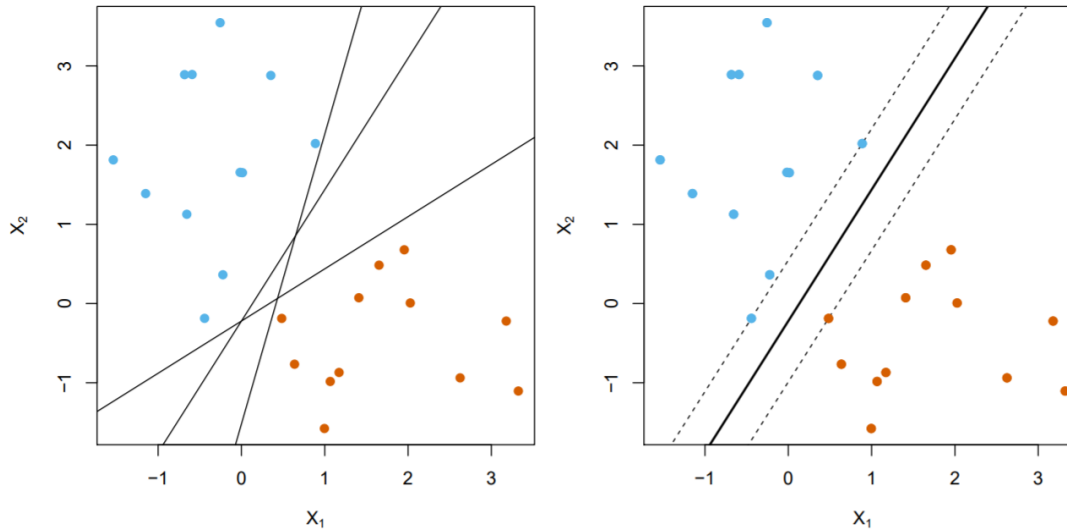
(in case of convergence)

Support Vector Machines

- /// Main Idea and Details
- /// Linear Support Vector Machines
- /// Kernel/Nonlinear Support Vector Machines
- /// Kernel selection and fine-tuning

Main Idea and Details

/// **Main idea:** given a $D = \{(x_i, y_i)\}_{i=1}^n$, $y_i \in \{-1, 1\}$, $H = \{x \in \mathbb{R}^J : \langle w, x \rangle = c\}$, with normal vector w and offset c , separating points from **class -1** and **class 1** with the largest symmetric **margin of $\pm\eta$**



/// **A contrast:** whilst Logistic Regression first build a conditional probability model, and then classify based in a threshold, SVMs bypass the first step, and build a classifier directly!

Linear Support Vector Machines

/// **Hard margin Linear SVM:** Find a hyperplane $H = \{x \in \mathbb{R}^J : \langle w, x \rangle = c\}$, with normal vector w and offset c , separating points from **class -1** and **class 1** with the largest symmetric **margin of** $\pm\eta = \frac{1}{||w||}$

$$\max_{w,c} \eta, \quad \text{s.t. } y_i \cdot (\langle w, x_i \rangle - c) \geq 1, \quad y_i \in \{-1, 1\}$$

Too strict in practice!

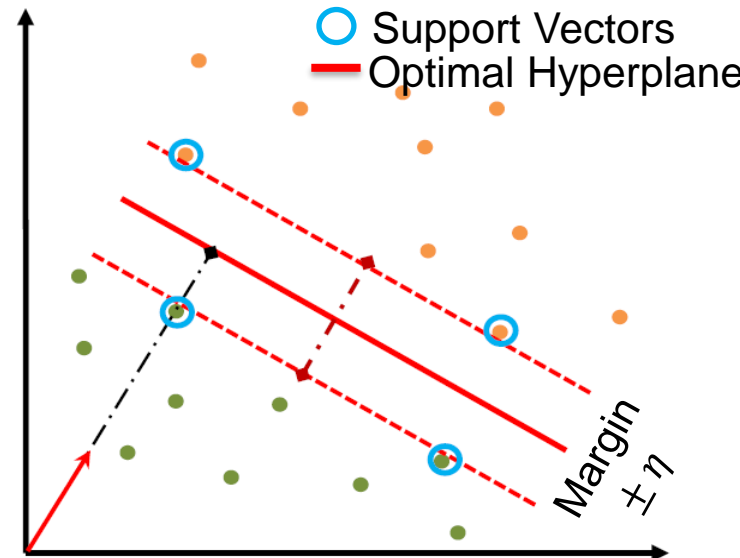
/// **Soft margin Linear SVM**

$$\min_{w,c} \frac{1}{2} ||w||^2 + C \sum_i \xi_i,$$

$$\text{s.t. } y_i \cdot (\langle w, x_i \rangle - c) \geq 1 - \xi_i \text{ and } \xi_i \geq 0$$

where ξ_i is a slack variable for each data point; it measures how much the i -th sample is allowed to violate the margin. Two conflicting objectives, maximize η but minimizing ξ_i as possible; C is the **hyperparameter that control this trade-off**

Both problems can be solved via Quadratic Programming solvers



Kernel Support Vector Machines

Primal form

$$\min_{\mathbf{w}, c} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i,$$

$$s.t. y_i \cdot (\langle \mathbf{w}, \mathbf{x}_i \rangle - c) \geq 1 - \xi_i \text{ and } \xi_i \geq 0$$

Dual form

$$\min_{\alpha} \frac{1}{2} \alpha^T \mathbf{Q} \alpha - \sum_i \alpha_i$$

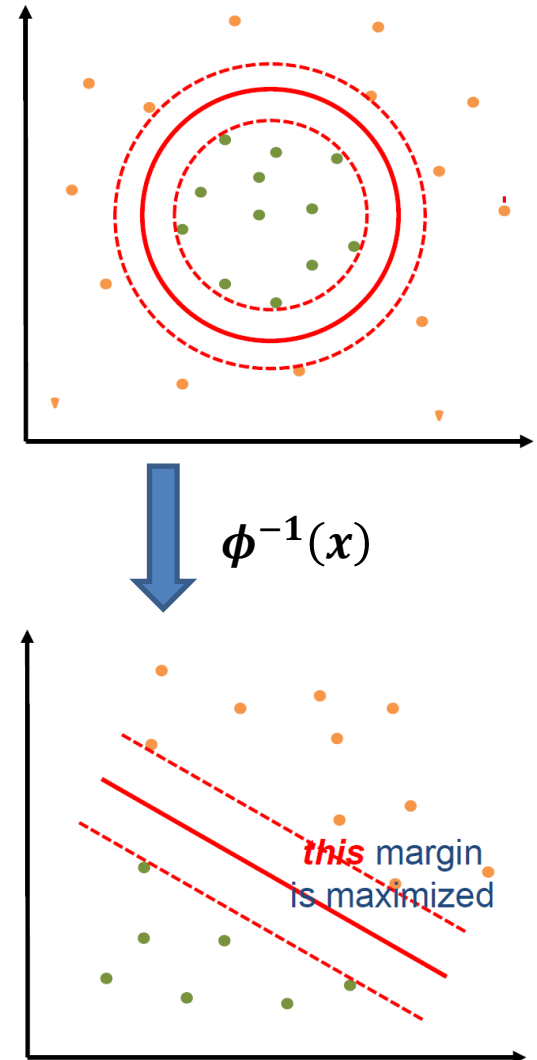
$$s.t. \mathbf{y}^T \alpha = 0, \text{ with } 0 \leq \alpha_i \leq C, i=1, \dots, n$$

with $\mathbf{Q}_{ik} := y_i y_k k(\mathbf{x}_i, \mathbf{x}_k)$ is an $n \times n$ positive definite matrix and C an upper bound for the dual weights

Decision function

$$f(\mathbf{z}) = \text{sgn} \left(\sum_i y_i \alpha_i K(\mathbf{x}_i, \mathbf{z}) + c \right)$$

Also solvable via QP, via Slater Theorem



Kernel selection and fine-tuning

- /// **Kernel Ridge Regression** (remember to scale input variables)
 - /// Start with a **Grid-search strategy, making it fine-grained** after a certain range looks promising
 - /// Try a range of regularization values: $\lambda = [0.0001, 1.0]$
 - /// RBF Kernel is preferred; values ranging: $\sigma = [0.01, 10.0]$
 - /// Use RMSE, **but also MAD and Pseudo-R2** (Square of the Correlation) for hyperparameter selection
- /// **Support Vector Machines** (remember to scale input variables)
 - /// **Grid-search strategy**: similar to KRR
 - /// Try a range of penalty values: $C = [0.001, 10.0]$
 - /// RBF Kernel is preferred; values ranging: $\sigma = [0.01, 10.0]$
 - /// **F1-Score tends to be a good generic metric**, but also Area Under the ROC curve (AUC), and precision/recall are important metrics. Accuracy is fine when the classes are well-balanced

Main Reading

- /// Chapters 5.8 and 12, of Elements of Statistical Learning, Hastie T., Tibshirani R. and Friedman J.
- /// Chapelle, O. (2007). Training a support vector machine in the primal. Neural computation, 19(5), 1155-1178.

email: faculty@mlinstitute.com

Further Reading

- /// Chapters 2 and 3 of Kernel Methods for Pattern Analysis, Shawe-Taylor J., and Cristianini N., Cambridge University Press
- /// Chapter 19, Support Vector Machines and Kernel Methods, Computer Age Statistical Inference, Efron B., Hastie T.
- /// Hofmann, T., Schölkopf, B., & Smola, A. J. (2008). Kernel methods in machine learning. The annals of statistics, 1171-1220.
- /// General Reference: Scholkopf, B., & Smola, A. J. (2001). Learning with kernels: support vector machines, regularization, optimization, and beyond. MIT press.

email: faculty@mliinstitute.com

DISCLAIMER

The material in this presentation is based on information that we consider reliable, but we do not warrant that it is accurate or complete, and it should not be relied on as such. Opinions expressed are current opinions only. We are not soliciting any action based upon this material. Neither the author, his employers, any operating arm of his employers nor any affiliated body can be held liable or responsible for any outcomes resulting from actions arising as a result of delivering this presentation. This presentation does not constitute investment advice nor should it be considered as such.

The views expressed in this presentation represent those of the lecturer in his or her individual private capacity and should not be taken to be the views of any employer or any affiliated body, including the MLI, or of the lecturer as an employee of any institution or affiliated body. Either he/she or his/her employers may or may not hold, or have recently held, a position in any security identified in this document.

This presentation is © MLI 2021. No part of this presentation may be copied, reproduced, distributed or stored in any form including electronically without express written permission in advance from the author.