Sapienza University of Rome

Master in Artificial Intelligence and Robotics
Master in Engineering in Computer Science

# Machine Learning

A.Y. 2020/2021

Prof. L. Iocchi, F. Patrizi

# 9. Kernel Methods

L. Iocchi, F. Patrizi

# Summary

- Kernel functions
- Kernelized linear models
- Kernelized SVM - classification
- Kernelized SVM - regression

*References*
C. Bishop. Pattern Recognition and Machine Learning. Chap. 6, Sect. 7.1

# Kernels

So far:
Objects represented as fixed-length feature-vectors $\mathbf{x} \in \mathbb{R}^M$ or $\phi(\mathbf{x})$.

Issue:
what about objects with variable length or infinite dimensions?

Examples:

- strings
- trees
- image features
- time-series
- ...

# Kernels

Approach:

use a *similarity measure* $k(\mathbf{x}, \mathbf{x}') \geq 0$ between the instances $\mathbf{x}, \mathbf{x}'$

$k(\mathbf{x}, \mathbf{x}')$ is called a *kernel* function.

Note: If we have $\phi(\mathbf{x})$ a possible choice is $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$.

for models which are based on a fixed nonlinear feature space mapping φ(x), the kernel function is given by this relation

# Kernels

**Definition**
*Kernel function*: a real-valued function $k(\mathbf{x}, \mathbf{x}') \in \mathbb{R}$, for $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$, where $\mathcal{X}$ is some abstract space.
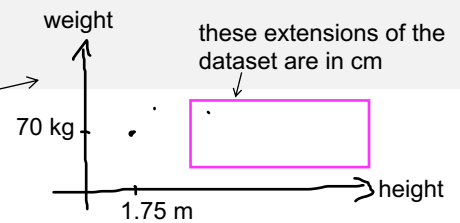
Typically $k$ is:
- symmetric: $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$
- non-negative: $k(\mathbf{x}, \mathbf{x}') \geq 0$.

Note: Not strictly required!

important is the normalization of data

# Input normalization

weight

these extensions of the dataset are in cm

70 kg

1.75 m

height

Input data in the dataset $D$ must be normalized in order for the kernel to be a good *similarity measure* in practice.

Several types of normalizations:

- min-max $\quad \bar{x} = \frac{x - min}{max - min}$    the min is 0 and the max is 1
  $min, max$: minimum and maximum input values in $D$
- normalization (standardization) $\quad \bar{x} = \frac{x - \mu}{\sigma}$   after this standardization the dataset will have mean=0 and standard deviation=1
  $\mu$ mean and $\sigma$ standard deviation of input values in $D$
- unit vector $\quad \bar{x} = \frac{x}{||x||}$

In the following, we assume the use of normalized input data.

# Kernel families

**Linear**

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$$

**Polynomial**

$$k(\mathbf{x}, \mathbf{x}') = (\beta \mathbf{x}^T \mathbf{x}' + \gamma)^d, \ d \in \{2, 3, \ldots\}$$

**Radial Basis Function (RBF)**

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\beta |\mathbf{x} - \mathbf{x}'|^2)$$

**Sigmoid**

$$k(\mathbf{x}, \mathbf{x}') = \tanh(\beta \mathbf{x}^T \mathbf{x}' + \gamma)$$

how do we use this kernel?

# Kernelized linear models

Consider a linear model $y(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}$ with dataset $D = \{(\mathbf{x}_n, t_n)_{n=1}^N\}$

we define an error function for ex. a squared error function + regularization term

Minimize $J(\mathbf{w}) = (\mathbf{t} - \mathbf{Xw})^T(\mathbf{t} - \mathbf{Xw}) + \lambda \|\mathbf{w}\|^2$

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} \text{ design matrix,} \qquad \mathbf{t} = \begin{bmatrix} t_1 \\ \vdots \\ t_N \end{bmatrix} \text{ output vector}$$

Optimal solution

we call it alpha

$\hat{\mathbf{w}} = (\mathbf{X}^T\mathbf{X} + \lambda I_N)^{-1}\mathbf{X}^T\mathbf{t} = \mathbf{X}^T\underline{(\mathbf{X}\mathbf{X}^T + \lambda I_N)^{-1}\mathbf{t}}$,
with $I_N$ the $N \times N$ identity matrix.    identity matrix NxN

# Kernelized linear models

Let $\boldsymbol{\alpha} = (\mathbf{X}\mathbf{X}^T + \lambda I_N)^{-1}\mathbf{t}$,

linear cobination of the samples in dataset x_n multiplied by coeff Alpha_n

then $\hat{\mathbf{w}} = \mathbf{X}^T\boldsymbol{\alpha} = \sum_{n=1}^N \alpha_n \mathbf{x}_n$.

Hence we have $y(\mathbf{x}; \hat{\mathbf{w}}) = \hat{\mathbf{w}}^T\mathbf{x} = \sum_{n=1}^N \alpha_n \mathbf{x}_n^T \mathbf{x}$.     x is used to indicate a general sample
while x_n is a specific sample

If we consider a linear kernel $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T\mathbf{x}'$, we can rewrite the model as

$$y(\mathbf{x}; \hat{\mathbf{w}}) = \sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x})$$

with

$\boldsymbol{\alpha} = (K + \lambda I_N)^{-1}\mathbf{t}$, and $K = \mathbf{X}\mathbf{X}^T$ **Gram matrix**

summarize

# Kernelized linear models

with linear kernel we will represent a linear function

## Linear model with linear kernel $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$

$$y(\mathbf{x}; \boldsymbol{\alpha}) = \sum_{n=1}^{N} \alpha_n \mathbf{x}_n^T \mathbf{x}$$

k(x_n, x)

## Solution

$$\boldsymbol{\alpha} = (K + \lambda I_N)^{-1} \mathbf{t}$$

## Gram matrix

should be positive semidefinite for all possible choices of the set {xn}

$$K = \begin{bmatrix} \mathbf{x}_1^T \mathbf{x}_1 & \cdots & \mathbf{x}_1^T \mathbf{x}_N \\ \vdots & \ddots & \vdots \\ \mathbf{x}_N^T \mathbf{x}_1 & \cdots & \mathbf{x}_N^T \mathbf{x}_N \end{bmatrix}$$

k ( x _ i , x _

cross product of all the possible pairs in the dataset

# Kernelized linear models  written with

## Linear model with any kernel $k$

$$y(\mathbf{x}; \boldsymbol{\alpha}) = \sum_{n=1}^{N} \alpha_n \, k(\mathbf{x}_n, \mathbf{x})$$

## Solution

$$\boldsymbol{\alpha} = (K + \lambda I_N)^{-1} \mathbf{t}$$

## Gram matrix

$$K = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

# Kernel trick

*Kernel trick* or *kernel substitution*

If input vector $\mathbf{x}$ appears in an algorithm only in the form of an inner product $\mathbf{x}^T \mathbf{x}'$, replace the inner product with some kernel $k(\mathbf{x}, \mathbf{x}')$.

- Can be applied to any $\mathbf{x}$ (even infinite size)
- No need to know $\phi(\mathbf{x})$
- Directly extend many well-known algorithms

now we examine the kernel-based algorithms that have sparse solutions, so that predictions for new inputs depend only on the kernel function evaluated at a subset of the training data points.

solving problems in classification, regression, and novelty detection.

# Kernelized SVM - classification

In SVM, solution has the form:

$$\hat{\mathbf{w}} = \sum_{n=1}^{N} \alpha_n \mathbf{x}_n$$

$$y(x) = w^T \phi(x) + b$$

Linear model (with linear kernel)

$$y(\mathbf{x}; \boldsymbol{\alpha}) = \text{sign}\left( w_0 + \sum_{n=1}^{N} \alpha_n \mathbf{x}_n^T \mathbf{x} \right)$$

Kernel trick

$$y(\mathbf{x}; \boldsymbol{\alpha}) = \text{sign}\left( w_0 + \sum_{n=1}^{N} \alpha_n \, k(\mathbf{x}_n, \mathbf{x}) \right)$$

Note: $w_0$ also estimated from $\boldsymbol{\alpha}$

# Kernelized SVM - classification

we have something like this:

$$x\_n^T x\_m$$

and we replace it with

Lagrangian problem for kernelized SVM classification

when the kernel is linear
we have the same result

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^{N} a_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} a_n \, a_m \, t_n \, t_m \, \boxed{k(\mathbf{x}_n, \mathbf{x}_m)}$$

**Solution**

$$a_n = \ldots$$

$$w_0 = \frac{1}{|SV|} \sum_{\mathbf{x}_i \in SV} \left( t_i - \sum_{\mathbf{x}_j \in S} a_j t_j k(\mathbf{x}_i, \mathbf{x}_j) \right)$$

# Kernelized linear regression

we have

Linear model for regression $y = \mathbf{w}^T \mathbf{x}$ and data set $D = \{(\mathbf{x}_n, t_n)_{n=1}^{N}\}$

Minimize the regularized loss function   error that measures how much the prediction is
distance from the value that we have in the dataset

$$J(\mathbf{w}) = \sum_{n=1}^{N} \boxed{E(y_n, t_n)} + \boxed{\lambda \|\mathbf{w}\|^2},$$

regularization term to control and avoid overfitting

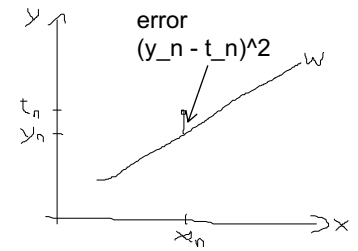where $y_n = \mathbf{w}^T \mathbf{x}_n$.

# Kernelized linear regression

Consider $E(y_n, t_n) = (y_n - t_n)^2$: i.e., regularized linear regression.

**Solution**

$$\boxed{\hat{\mathbf{w}} = (\mathbf{X}^T\mathbf{X} + \lambda I_N)^{-1}\mathbf{X}^T\mathbf{t} = \mathbf{X}^T\boldsymbol{\alpha}}$$

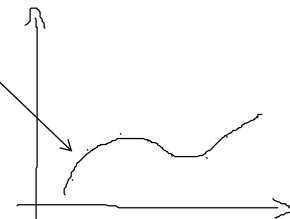$$\boldsymbol{\alpha} = (\mathbf{X}\mathbf{X}^T + \lambda I_N)^{-1}\mathbf{t}$$

error
(y_n - t_n)^2

Predictions are made using:

we can apply the kernel trick and
so we have the kernelized model

$$y(\mathbf{x}; \hat{\mathbf{w}}) = \sum_{n=1}^{N} \alpha_n \boxed{\mathbf{x}_n^T \mathbf{x}.}$$

# Kernelized linear regression

for ex. if we apply a polynomial kernel
(cubic with 3 degrees)

Apply the kernel trick:

$$y(\mathbf{x}; \hat{\mathbf{w}}) = \sum_{n=1}^{N} \alpha_n k(\mathbf{x}_n, \mathbf{x})$$

$$\boldsymbol{\alpha} = (K + \lambda I_N)^{-1}\mathbf{t}$$

Issue: computation of $K$ requires $|D|^2$ operations and $K$ is not sparse.

this method may not be practical when the dataset is large

# Kernelized SVM - regression

Consider

we have a different error function, it is not a quadratic functio

$$J(\mathbf{w}) = C \sum_{n=1}^{N} E_\epsilon(y_n, t_n) + \frac{1}{2}\|\mathbf{w}\|^2,$$
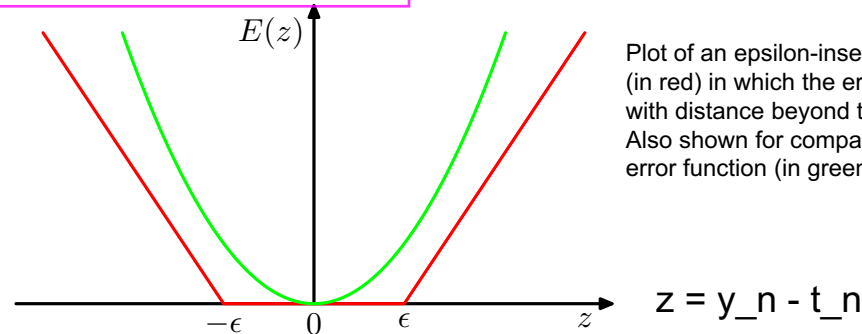
inverse of normalization term

with $C$ inverse of $\lambda$ and an $\epsilon$-insensitive error function:

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

← if the distance between t and y is less then epsilon

error is proportional to the distance from the epsilon boundary

Plot of an epsilon-insensitive error function (in red) in which the error increases linearly with distance beyond the insensitive region. Also shown for comparison is the quadratic error function (in green).
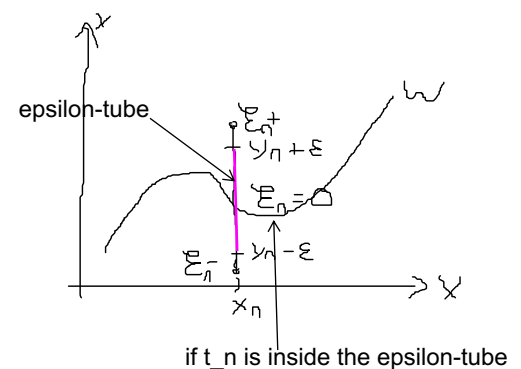
z = y_n - t_n

Not differentiable $\rightarrow$ difficult to solve.

# Kernelized SVM - regression

we can re-express the optimization problem by introducing slack variables.

Introduce *slack variables* $\xi_n^+, \xi_n^- \geq 0$:

$$t_n \leq y_n + \epsilon + \xi_n^+$$
$$t_n \geq y_n - \epsilon - \xi_n^-$$
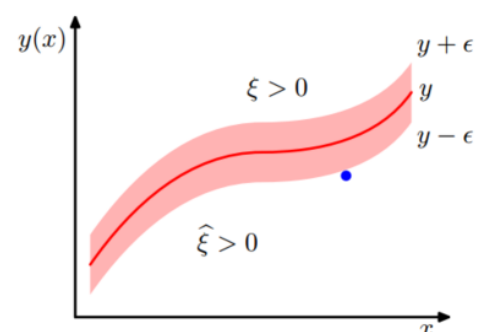
epsilon-tube

if t_n is inside the epsilon-tube

Points inside the $\epsilon$-tube $y_n - \epsilon \leq t_n \leq y_n + \epsilon \Rightarrow \xi_n = 0$

$\xi_n^+ > 0 \Rightarrow t_n > y_n + \epsilon$

$\xi_n^- > 0 \Rightarrow t_n < y_n - \epsilon$

with $y_n = y(\mathbf{x}_n; \mathbf{w})$

Illustration of SVM regression, showing the regression curve together with the $\epsilon$-insensitive 'tube'. Also shown are examples of the slack variables $\xi$ and $\hat{\xi}$. Points above the $\epsilon$-tube have $\xi > 0$ and $\hat{\xi} = 0$, points below the $\epsilon$-tube have $\xi = 0$ and $\hat{\xi} > 0$, and points inside the $\epsilon$-tube have $\xi = \hat{\xi} = 0$.

# Kernelized SVM - regression

epsilon is up to we, we need to choose a good value for it

Loss function can be rewritten as:

$$J(\mathbf{w}) = C \sum_{n=1}^{N}(\xi_n^+ + \xi_n^-) + \frac{1}{2}\|\mathbf{w}\|^2,$$

subject to the constraints:

$$t_n \leq y(\mathbf{x}_n; \mathbf{w}) + \epsilon + \xi_n^+$$
$$t_n \geq y(\mathbf{x}_n; \mathbf{w}) - \epsilon - \xi_n^-$$
$$\xi_n^+ \geq 0$$
$$\xi_n^- \geq 0$$

This is a standard quadratic program (QP), can be "easily" solved.

# Kernelized SVM - regression

Lagrangian problem

$$\tilde{L}(\mathbf{a}, \mathbf{a}') = \ldots \sum_{n=1}^{N}\sum_{m=1}^{N} a_n a_m \ldots k(\mathbf{x}_n, \mathbf{x}_m) \ldots$$

from which we compute $\hat{a}_n$, $\hat{a}'_m$ (sparse values, most of them are zero) and

$$\hat{w}_0 = t_n - \epsilon - \sum_{m=1}^{N}(\hat{a}_m - \hat{a}'_m)\, k(\mathbf{x}_n, \mathbf{x}_m)$$

for some data point $n$ such that $0 < a_n < C$

**Prediction**

$$y(\mathbf{x}) = \sum_{n=1}^{N}(\hat{a}_n - \hat{a}'_n)\, k(\mathbf{x}, \mathbf{x}_n) + \hat{w}_0$$

# Kernelized SVM - regression

one of the best solution to solve ML problems

From Karush-Kuhn-Tucker (KKT) condition (see Bishop Sect. 7.1.4)
**Support vectors** contribute to predictions

$\hat{a}_n > 0 \Rightarrow \epsilon + \xi_n + y_n - t_n = 0$   for all the points inside the epsilon-tube that not contribute to the solution, the lagrangian is zero. For these points we don't need to compute the kernel
data point lies on or above upper boundary of the $\epsilon$-tube

$\hat{a}'_n > 0 \Rightarrow \epsilon + \xi_n - y_n + t_n = 0$
data point lies on or below lower boundary of the $\epsilon$-tube

All other data points inside the $\epsilon$-tube have $\hat{a}_n = 0$ and $\hat{a}'_n = 0$ and thus do not contribute to prediction.

# Kernelized SVM - regression

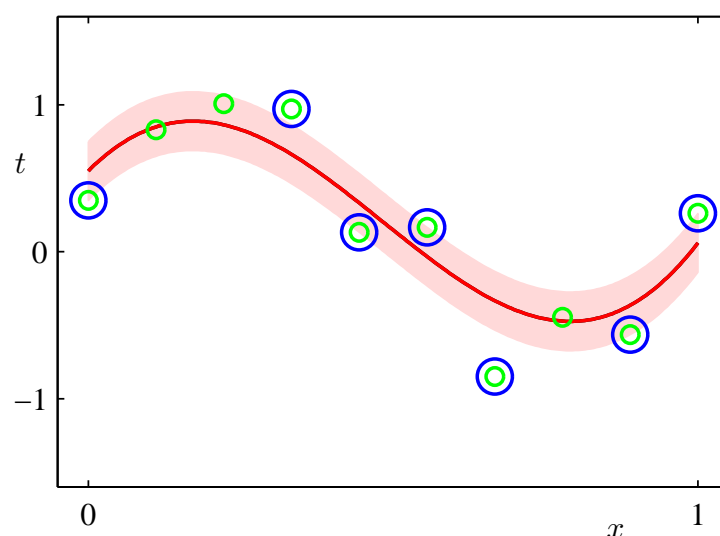Example: support vectors and $\epsilon$ insensitive tube



Illustration of the v-SVM for regression applied to the sinusoidal synthetic data set using Gaussian kernels. The predicted regression curve is shown by the red line, and the -insensitive tube corresponds to the shaded region. Also, the data points are shown in green, and those with support vectors are indicated by blue circles.

# Summary

- Kernel methods overcome difficulties in defining non-linear models
- Kernelized SVM is one of the most effective ML method for classification and regression
- Still requires model selection and hyper-parameters tuning