Introduction to Kernel Methods

Fabio A. González Ph.D.

The Kernel Approach to Machine Learning

The Kernel Trick

A Kernel Pattern Analysis Algorithm

Kernel Functions

Kernel Algorithms

Kernels in Complex Structured Data

# Introduction to Kernel Methods

Fabio A. González Ph.D.

Depto. de Ing. de Sistemas e Industrial
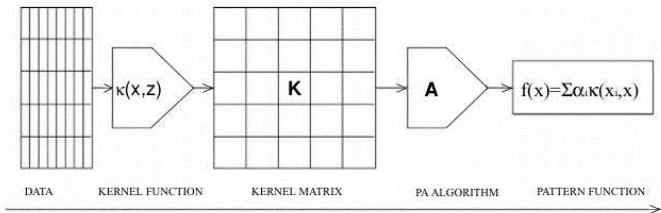Universidad Nacional de Colombia, Bogotá

March 13, 2007

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Outline

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

Summary

A modular process
for machine learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# The Approach

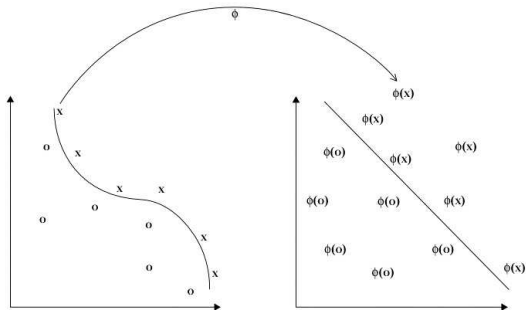- Data items are embedded into a vector space called the feature space
- Linear relations are sought among the images of the data items in the feature space
- The pattern analysis algorithm are based only on the pairwise dot products, they do not need the actual coordinates of the embedded points
- The pairwise dot products in the feature space could be efficiently calculated using a kernel function

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning
Summary
A modular process
for machine learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# The Process



DATA · KERNEL FUNCTION $\kappa(x,z)$ · KERNEL MATRIX $\mathbf{K}$ · PA ALGORITHM $\mathbf{A}$ · PATTERN FUNCTION $f(x)=\Sigma\alpha_i\kappa(x_i,x)$

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

Mapping the input
space to the feature
space

Calculating the dot
product in the
feature space

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Input space vs. feature space

- Why do we want to map to a different feature space?

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

Mapping the input
space to the feature
space

Calculating the dot
product in the
feature space

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Example (1)

- How to separate these two classes?

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

Mapping the input
space to the feature
space

Calculating the dot
product in the
feature space

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

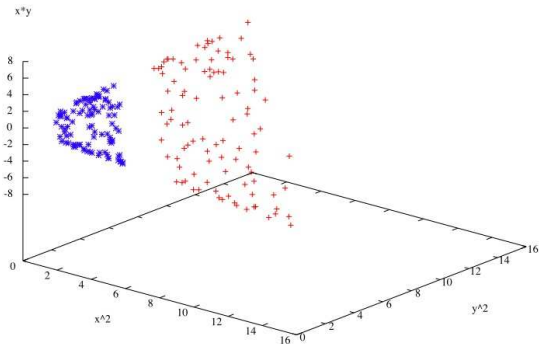# Example (2)

- Map to $\mathbb{R}^3$:

$$\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$
$$(x, y) \mapsto (x^2, y^2, xy)$$

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

Mapping the input
space to the feature
space

Calculating the dot
product in the
feature space

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Example (3)

• Map to $\mathbb{R}^3$:

$$\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$
$$(x, y) \mapsto (x^2, y^2, xy)$$

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

Mapping the input
space to the feature
space

Calculating the dot
product in the
feature space

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Dot product in the feature space

- 

$$\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$
$$(x_1, x_2) \mapsto (x_1^2, x_2^2, \sqrt{2}x_1 x_2)$$

- 

$$
\begin{aligned}
\langle \phi(x), \phi(z) \rangle &= \left\langle (x_1^2, x_2^2, \sqrt{2}x_1 x_2), (z_1^2, z_2^2, \sqrt{2}z_1 z_2) \right\rangle \\
&= x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 x_2 z_1 z_2 \\
&= (x_1 z_1 + x_2 z_2)^2 \\
&= \langle x, z \rangle^2
\end{aligned}
$$

- A function $k : X \times X \rightarrow \mathbb{R}$ such that
  $k(x, z) = \langle \phi(x), \phi(z) \rangle$ is called a kernel

- <u>Morale</u>: **you don't need to apply $\phi$ explicitly to calculate the dot product in the feature space!**

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

Mapping the input
space to the feature
space

Calculating the dot
product in the
feature space

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Kernel induced feature space

- The feature space induced by the kernel is not unique:
  The kernel

  $$k(x, z) = \langle x, z \rangle^2$$

  also calculates the dot product in the four dimensional
  feature space:

  $$\phi : \mathbb{R}^2 \quad \rightarrow \quad \mathbb{R}^4$$
  $$(x_1, x_2) \quad \mapsto \quad (x_1^2, x_2^2, x_1 x_2, x_2 x_1)$$

- The example can be generalised to $\mathbb{R}^n$

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Primal linear
regression

Dual linear regression

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Problem definition

- Given a training set $S = \{(x_1, y_1), \ldots, (x_l, y_l)\}$ of points $x_i \in \mathbb{R}^n$ with corresponding labels $y_i \in \mathbb{R}$ the problem is to find a real-valued linear function that best interpolates the training set:

$$g(x) = \langle w, x \rangle = w'x = \sum_{i=1}^{n} w_i x_i$$

- If the data points were generated by a function like $g(x)$, it is possible to find the parameters $w$ by solving

$$Xw = y$$

where

$$X = \begin{bmatrix} x'_1 \\ \vdots \\ x'_l \end{bmatrix}$$

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Primal linear
regression

Dual linear regression

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Graphical representation

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Primal linear
regression

Dual linear regression

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Loss function

- Minimize

$$\mathcal{L}(g, S) = \mathcal{L}(\mathrm{w}, S) = \sum_{i=1}^{l} (y_i - g(x_i))^2 = \sum_{i=1}^{l} \xi_i^2$$

$$= \sum_{i=1}^{l} \mathcal{L}(g, (\mathrm{x}_i, y_i))$$

- This could be written as

$$\mathcal{L}(\mathrm{w}, S) = \|\xi\|^2 = (\mathrm{y} - \mathrm{Xw})'(\mathrm{y} - \mathrm{Xw})$$

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Primal linear
regression

Dual linear regression

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Solution

$$\frac{\partial \mathcal{L}(\mathrm{w}, S)}{\partial \mathrm{w}} = -2\mathrm{X}'\mathrm{y} + 2\mathrm{X}'\mathrm{Xw} = 0,$$

therefore

$$\mathrm{X}'\mathrm{Xw} = \mathrm{X}'\mathrm{y},$$

and

$$\mathrm{w} = (\mathrm{X}'\mathrm{X})^{-1}\mathrm{X}'\mathrm{y}$$

# Dual representation of the problem

$$\mathrm{w} = (\mathrm{X'X})^{-1}\mathrm{X'y} = \mathrm{X'X}(\mathrm{X'X})^{-2}\mathrm{X'y} = \mathrm{X'}\alpha$$

- So, $\mathrm{w}$ is a linear combination of the training samples,
  $\mathrm{w} = \sum_{i=1}^{l} \alpha_i \mathrm{x}_i$.

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Primal linear
regression

Dual linear regression

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Solution

- From the solution of the primal problem:

$$X'Xw = X'y,$$

- then

$$XX'Xw = XX'y,$$

- using the dual representation

$$XX'XX'\alpha = XX'y,$$

- then

$$\alpha = (XX')^{-1}y,$$

- and

$$g(x) = w'x = \alpha'Xx.$$

- <u>Note</u>: $XX'$ may be close to singular, or singular according to machine precision.

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Primal linear
regression

Dual linear regression

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Ridge regression

- If $XX'$ is singular, the pseudo-inverse could be used: to find the $w$ that satisfies $X'Xw = X'y$ with minimal norm.

- Optimisation problem:

$$\min_{w} \mathcal{L}_\lambda(w, S) = \min_{w} \lambda \|w\|^2 + \sum_{i=1}^{l} (y_i - g(x_i))^2,$$

where $\lambda$ defines the trade-off between norm and loss. This controls the complexity of the model (the porcess is called *regularization*).

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Primal linear
regression

Dual linear regression

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

## Solution

- Taking the derivative and making it equal to zero:

$$X'Xw + \lambda w = (X'X + \lambda I_n w) = X'y,$$

- then,

$$w = (X'X + \lambda I_n)^{-1} X'y.$$

- In terms of $\alpha$:

$$w = \lambda^{-1} X'(y - Xw) = X'\alpha,$$

- then

$$\alpha = \lambda^{-1}(y - Xw) = (XX' + \lambda I_l)^{-1}y.$$

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Primal linear
regression

Dual linear regression

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

## Prediction function

$$\begin{aligned} g(\mathrm{x}) = \langle \mathrm{w}, \mathrm{x} \rangle &= \left\langle \sum_{i=1}^{l} \alpha_i \mathrm{x_i}, \mathrm{x} \right\rangle = \sum_{i=1}^{l} \alpha_i \langle \mathrm{x_i}, \mathrm{x} \rangle \\ &= y'(\mathrm{G} + \lambda \mathrm{I}_l)^{-1} \mathrm{k}, \end{aligned}$$

where $\mathrm{G} = \mathrm{XX}'$ (called the Gram Matrix) and $\mathrm{k}_i = \langle \mathrm{x_i}, \mathrm{x} \rangle$.

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Mathematical
characterisation

Visualizing kernels in
input space

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Characterisation

### Theorem
*A function*

$$k : X \times X \to \mathbb{R},$$

*which is either continuous or has a countable domain, can be decomposed*

$$k(\mathrm{x}, \mathrm{z}) = \langle \phi(\mathrm{x}), \phi(\mathrm{z}) \rangle$$

*into a feature map $\phi$ into a Hilbert space $F$ applied to both its arguments followed by the evaluation of the inner product in $F$ if and only if it satisfies the finitely positive semi-definite property.*

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Mathematical
characterisation

Visualizing kernels in
input space

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Some kernel functions

Assume $k_1$ and $k_2$ kernels:

- $k(\mathrm{x}, \mathrm{z}) = p(k_1(\mathrm{x}, \mathrm{z}))$. $p$ a polynomial with positive coefficients.

- $k(\mathrm{x}, \mathrm{z}) = \exp(k_1(\mathrm{x}, \mathrm{z}))$.

- $k(\mathrm{x}, \mathrm{z}) = \exp(-\|\mathrm{x} - \mathrm{z}\|^2 / (2\sigma^2))$. Gaussian kernel.

- $k(\mathrm{x}, \mathrm{z}) = k_1(\mathrm{x}, \mathrm{z})k_2(\mathrm{x}, \mathrm{z})$

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Mathematical
characterisation

Visualizing kernels in
input space

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Embeddings corresponding to kernels

- It is possible to calculate the feature space induced by a kernel (Mercer's Theorem)
- This can be done in a constructive way
- The feature space can even be of infinite dimension.

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Mathematical
characterisation

Visualizing kernels in
input space

Kernel
Algorithms

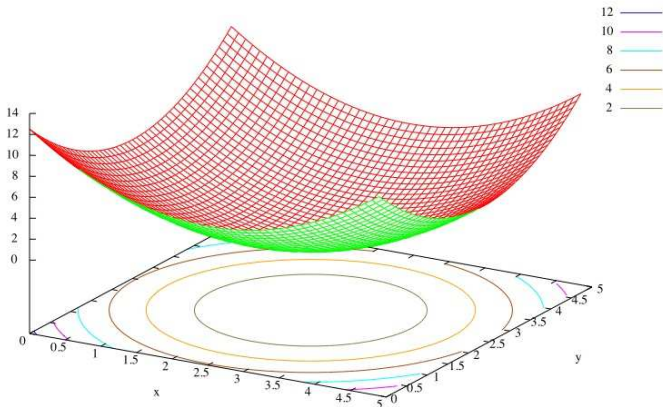Kernels in
Complex
Structured
Data

# How to visualize?

- Choose a point in input space $p_0$
- Calculate the distance from another point $x$ to $p_0$ in the feature space:

$$
\begin{aligned}
\|\phi(p_0) - \phi(x)\|_F^2 &= \langle \phi(p_0) - \phi(x), \phi(p_0) - \phi(x) \rangle_F \\
&= \langle \phi(p_0), \phi(p_0) \rangle_F + \langle \phi(x), \phi(x) \rangle_F \\
&\quad - 2 \langle \phi(p_0), \phi(x) \rangle_F \\
&= k(p_0, p_0) + k(x, x) - 2k(p_0, x)
\end{aligned}
$$

- Plot $f(x) = \|\phi(p_0) - \phi(x)\|_F^2$

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions
Mathematical
characterisation
Visualizing kernels in
input space

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Identity kernel

$$k(x, z) = \langle x, z \rangle$$

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions
Mathematical
characterisation
Visualizing kernels in
input space

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Quadratic kernel (1)

$$k(x, z) = \langle x, z \rangle^2$$

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions
Mathematical
characterisation
Visualizing kernels in
input space

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Identity kernel (2)

$$k(x, z) = \langle x, z \rangle^2$$

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Mathematical
characterisation

Visualizing kernels in
input space

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Gaussian kernel

$$k(\mathrm{x}, \mathrm{z}) = e^{-\frac{\|\mathrm{x}-\mathrm{z}\|^2}{2\sigma^2}}$$

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Basic computations in feature space

- Means
- Distances
- Projections
- Covariance

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Classification and regression

- Support Vector Machines
- Support Vector Regression
- Kernel Fisher Discriminant
- Kernel Perceptron

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Dimensionality reduction and clustering

- Kernel PCA
- Kernel CCA
- Kernel $k$-means
- Kernel SOM

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Kernels in complex structured data

- Since kernel methods do not require an attribute-based representation of objects, it is possible to perform learning over complex structured data (or unstructured data)
- We only need to define a dot product operation (similarity, dissimilarity measure)
- Examples:
  - Strings
  - Texts
  - Trees
  - Graphs

Introduction
to Kernel
Methods

Fabio A.
González
Ph.D.

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# References

📄 Shawe-Taylor, J. and Cristianini, N. 2004 Kernel Methods for Pattern Analysis. Cambridge University Press.