# Introduction to Kernel Methods

## Fabio A. González Ph.D.

Depto. de Ing. de Sistemas e Industrial
Universidad Nacional de Colombia, Bogotá
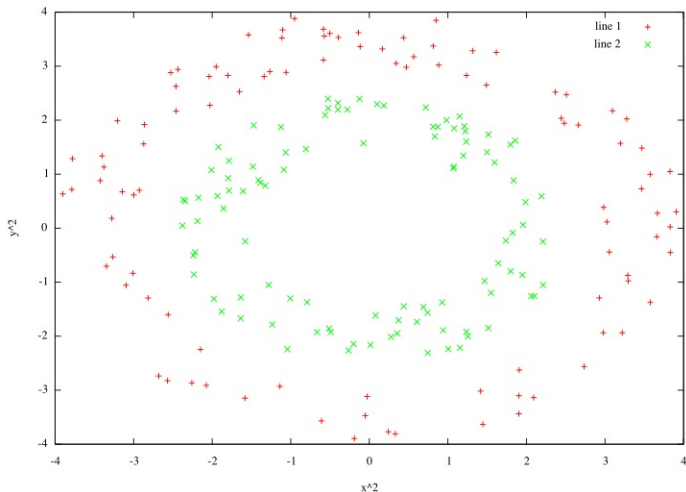
March 13, 2008

Introduction
to Kernel
Methods

F. González

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

F. González

The Kernel
Approach to
Machine
Learning
Motivation
Overview

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Problem 1

How to separate these two classes using a linear function?

Introduction
to Kernel
Methods

F. González

The Kernel
Approach to
Machine
Learning
Motivation
Overview

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
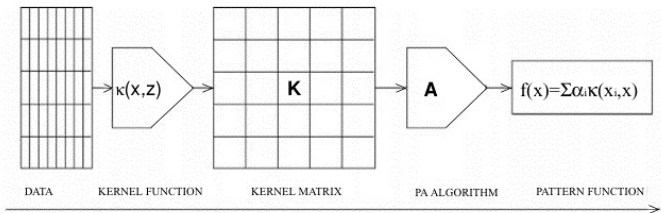Algorithms

Kernels in
Complex
Structured
Data

# Problem 2

How to do symbolic regression?

$$\Sigma = \{A, C, G, T\}$$

$$
\begin{array}{cccc}
f: & \Sigma^d & \to & \mathbb{R} \\
 & ACGTA & \mapsto & 10.0 \\
 & GTCCA & \mapsto & 11.3 \\
 & GGTAC & \mapsto & 1.0 \\
 & CCTGA & \mapsto & 4.5 \\
 & \vdots & \vdots & \vdots
\end{array}
$$

Introduction
to Kernel
Methods

F. González

The Kernel
Approach to
Machine
Learning
Motivation
Overview

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# The Process



DATA    KERNEL FUNCTION    KERNEL MATRIX    PA ALGORITHM    PATTERN FUNCTION

Introduction
to Kernel
Methods

F. González

The Kernel
Approach to
Machine
Learning
Motivation
Overview

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

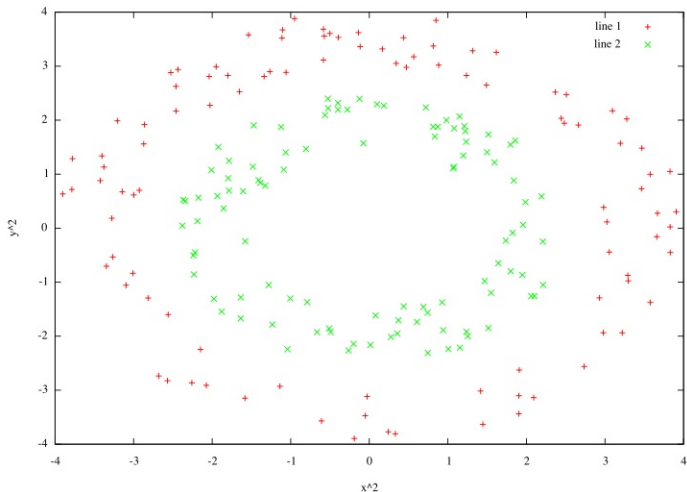Kernels in
Complex
Structured
Data

# 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
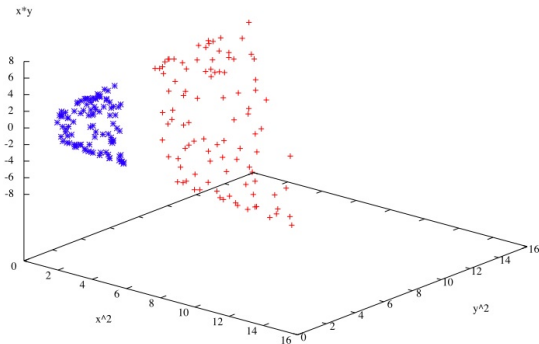- The pairwise dot products in the feature space could be efficiently calculated using a kernel function

Introduction
to Kernel
Methods

F. González

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

# Problem 1

- How to separate these two classes using a linear function?

# Solution

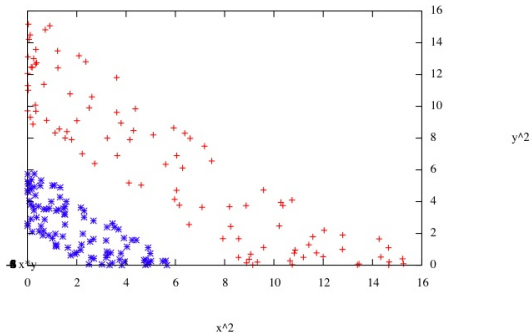- Map to $\mathbb{R}^3$:

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

Introduction
to Kernel
Methods

F. González

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

# Solution

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

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



x^2

Introduction
to Kernel
Methods

F. González

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

Introduction
to Kernel
Methods

F. González

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

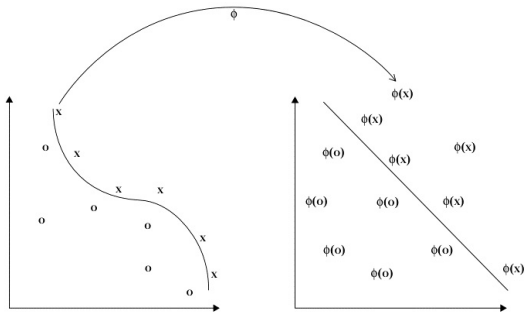Mapping the input
space to the feature
space

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

- 

$$\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$
$$(x_1, x_2) \mapsto (x_1^2, x_2^2, \sqrt{2}x_1x_2)$$

- 

$$
\begin{aligned}
\langle \phi(x), \phi(z) \rangle &= \left\langle (x_1^2, x_2^2, \sqrt{2}x_1x_2), (z_1^2, z_2^2, \sqrt{2}z_1z_2) \right\rangle \\
&= x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1x_2z_1z_2 \\
&= (x_1z_1 + x_2z_2)^2 \\
&= \langle x, z \rangle^2
\end{aligned}
$$

- 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!**

# 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:

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

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

Introduction
to Kernel
Methods

F. González

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

F. González

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

F. González

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

F. González

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

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

therefore

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

and

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

Introduction
to Kernel
Methods

F. González

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

# 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

F. González

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

F. González

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 process is called *regularization*).

# Solution

- Taking the derivative and making it equal to zero:

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

where $I_n$ is an identity matrix of $n \times n$ dimension,

- 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

F. González

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

$$g(\mathrm{x}) = \langle \mathrm{w}, \mathrm{x} \rangle \;\; = \;\; \left\langle \sum_{i=1}^{l} \alpha_i \mathrm{x_i}, \mathrm{x} \right\rangle = \sum_{i=1}^{l} \alpha_i \langle \mathrm{x_i}, \mathrm{x} \rangle$$

Introduction
to Kernel
Methods

F. González

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 as a kernel method

- The Gram matrix $G = XX'$ is the matrix of dot products

$$
G = XX' = \begin{bmatrix} x'_1 \\ \vdots \\ x'_l \end{bmatrix} [x_1 \cdots x_l] = \begin{bmatrix} \langle x_1, x_1 \rangle & \langle x_1, x_l \rangle \\ \\ \langle x_l, x_1 \rangle & \langle x_l, x_l \rangle \end{bmatrix}
$$

- $G$ may be replaced by a general kernel matrix, $K$, with $k_{ij} = k(x_i, x_j) = <\phi(x_i), \phi(x_j)>$
- The $\alpha$'s are calculated as:

$$
\alpha = (K + \lambda I_l)^{-1} y
$$

- The predicted function is approximated as:

$$
g(x) = \sum_{i=1}^{l} \alpha_i k(x, x_i) = y'(K + \lambda I_l)^{-1} \begin{bmatrix} k(x, x_1) \\ \vdots \\ k(x, x_l) \end{bmatrix}
$$

Introduction
to Kernel
Methods

F. González

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
*(Mercer's 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

F. González

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

F. González

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
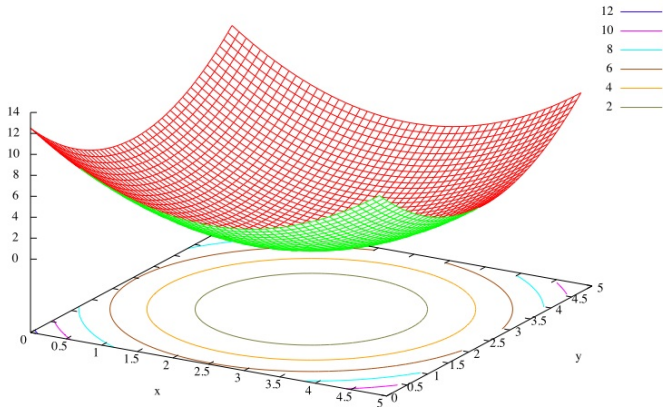- The feature space can even be of infinite dimension.

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

# Identity kernel

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

Introduction
to Kernel
Methods

F. González

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

F. González

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

F. González

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

F. González

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

F. González

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

F. González

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

F. González

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

F. González

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Problem 2

How to do symbolic regression?

$$\Sigma = \{A, C, G, T\}$$

$$
\begin{array}{rccc}
f: & \Sigma^d & \rightarrow & \mathbb{R} \\
& ACGTA & \mapsto & 10.0 \\
& GTCCA & \mapsto & 11.3 \\
& GGTAC & \mapsto & 1.0 \\
& CCTGA & \mapsto & 4.5 \\
& \vdots & \vdots & \vdots
\end{array}
$$

Introduction
to Kernel
Methods

F. González

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Solution

- Define a kernel on strings

$$k : \Sigma^d \times \Sigma^d \to \mathbb{R}$$

- Use the kernel along with a kernel learning regression algorithm to find the regression function
- What is a good candidate for $k$?
  - a function that measures string similarity
  - higher value for similar strings, smaller value for different strings
- 

$$k(s_1 \ldots s_d, t_1 \ldots t_d) = \sum_{i=1}^{n} equal(s_i, t_i)$$

$$equal(s_i, t_i) = \begin{cases} 1 & \text{if } s_i = t_i \\ 0 & \text{otherwise} \end{cases}$$

- $k(ACTAG, CCTCG) = ?$
- Is it a kernel?

Introduction
to Kernel
Methods

F. González

The Kernel
Approach to
Machine
Learning

The Kernel
Trick

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

# Induced Feature Space

- What is the feature space induced by $k$?
- 

$$\phi : \Sigma^d \;\; \to \;\; \mathbb{R}^{4d}$$
$$s_1 \dots s_d \;\; \mapsto \;\; (x_1^1, \dots, x_4^1, x_1^2, \dots, x_4^2, \dots, x_1^d, \dots, x_4^d)$$

$$(x_1^j, \dots, x_4^j) = \begin{cases} (1, 0, 0, 0) & \text{if } s_j = {'A'} \\ (0, 1, 0, 0) & \text{if } s_j = {'C'} \\ (0, 0, 1, 0) & \text{if } s_j = {'G'} \\ (0, 0, 0, 1) & \text{if } s_j = {'T'} \end{cases}$$

Introduction to Kernel Methods

F. González

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.