

Introduction to Kernel Methods

Fabio A. González Ph.D.

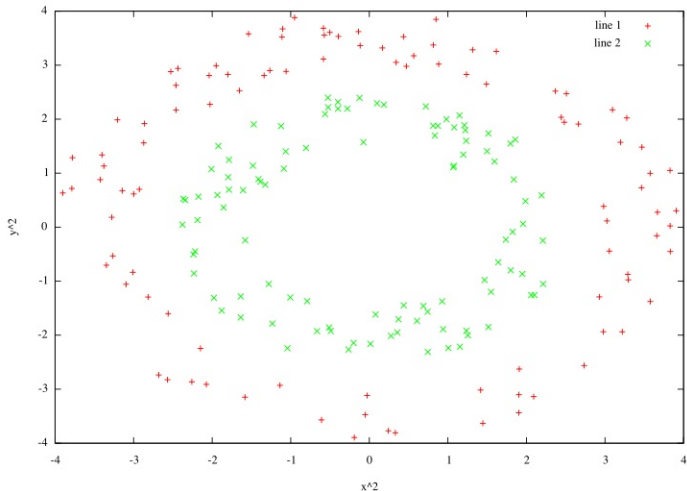
Depto. de Ing. de Sistemas e Industrial
Universidad Nacional de Colombia, Bogotá

April 2, 2009

- 1 Introduction
Motivation
- 2 The Kernel Trick
Mapping the input space to the feature space
Calculating the dot product in the feature space
- 3 The Kernel Approach to Machine Learning
- 4 A Kernel Pattern Analysis Algorithm
Primal linear regression
Dual linear regression
- 5 Kernel Functions
Mathematical characterisation
Visualizing kernels in input space
- 6 Kernel Algorithms
- 7 Kernels in Complex Structured Data

Problem 1

How to separate these two classes using a linear function?



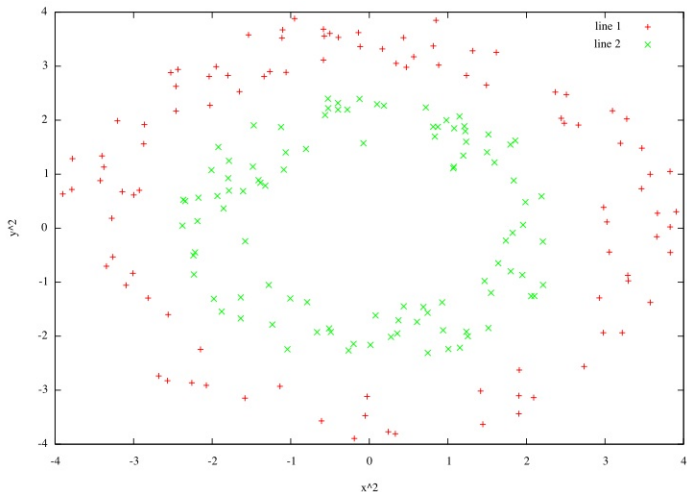
How to do symbolic regression?

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

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

Problem 1

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



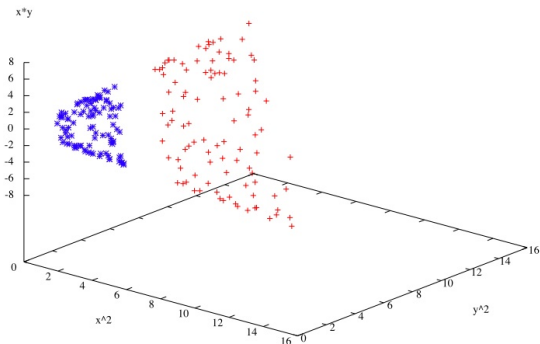
Mapping the input
space to the feature
space

Calculating the dot
product in the
feature space

- Map to \mathbb{R}^3 :

$$\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$(x, y) \mapsto (x^2, y^2, xy)$$

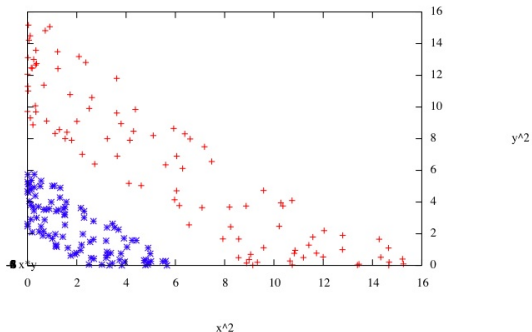


Mapping the input
space to the feature
space

Calculating the dot
product in the
feature space

- Map to \mathbb{R}^3 :

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



Input space vs. feature space

Introduction

The Kernel Trick

Mapping the input
space to the feature
space

Calculating the dot
product in the
feature space

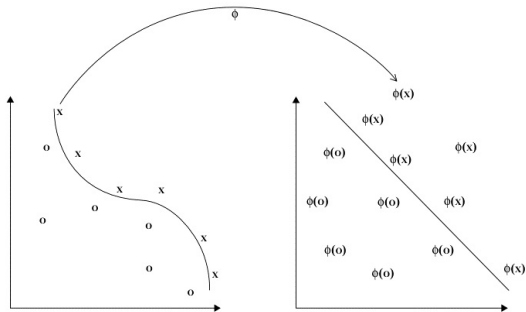
The Kernel Approach to Machine Learning

A Kernel Pattern Analysis Algorithm

Kernel Functions

Kernel Algorithms

Kernels in Complex Structured Data



Dot product in the feature space

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

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

- A function $k : X \times X \rightarrow \mathbb{R}$ such that $k(x, z) = \langle \phi(x), \phi(z) \rangle$ is called a kernel
- **Morale:** you don't need to apply ϕ explicitly to calculate the dot product in the feature space!

Kernel induced feature space

Introduction

The Kernel
Trick

Mapping the input
space to the feature
space

Calculating the dot
product in the
feature space

The Kernel
Approach to
Machine
Learning

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

- 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:

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

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

The Process

Introduction

The Kernel
Trick

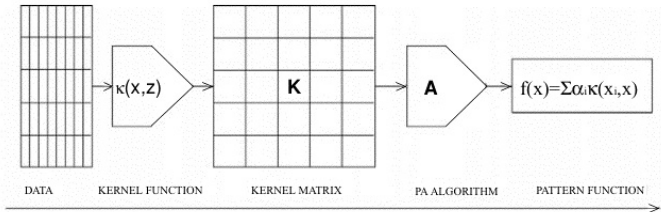
The Kernel
Approach to
Machine
Learning

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 algorithms 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

Problem definition

- Given a training set $S = \{(x_1, y_1), \dots, (x_l, y_l)\}$ 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:

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

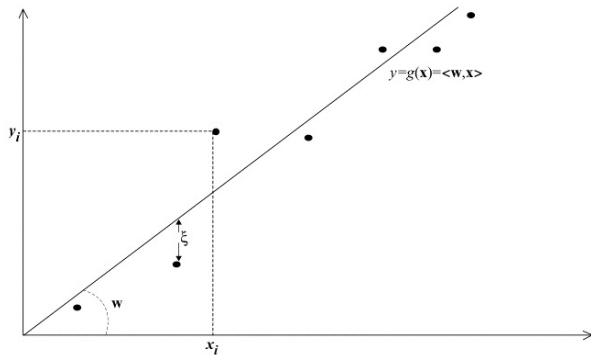
- If the data points were generated by a function like $g(x)$, it is possible to find the parameters w by solving

$$Xw = y$$

where

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

Graphical representation



Introduction

The Kernel
Trick

The Kernel
Approach to
Machine
Learning

A Kernel
Pattern
Analysis
Algorithm

Primal linear
regression

Dual linear regression

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

Loss function

- Minimize

$$\begin{aligned}\mathcal{L}(g, S) &= \mathcal{L}(\mathbf{w}, S) = \sum_{i=1}^l (y_i - g(x_i))^2 = \sum_{i=1}^l \xi_i^2 \\ &= \sum_{i=1}^l \mathcal{L}(g, (x_i, y_i))\end{aligned}$$

- This could be written as

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

Introduction

The Kernel
Trick

The Kernel
Approach to
Machine
Learning

A Kernel
Pattern
Analysis
Algorithm

Primal linear
regression

Dual linear regression

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

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

therefore

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

and

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

Dual representation of the problem

Introduction

The Kernel
Trick

The Kernel
Approach to
Machine
Learning

A Kernel
Pattern
Analysis
Algorithm

Primal linear
regression

Dual linear regression

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

- $\mathbf{w} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-2}\mathbf{X}'\mathbf{y} = \mathbf{X}'\boldsymbol{\alpha}$
- So, \mathbf{w} is a linear combination of the training samples,
$$\mathbf{w} = \sum_{i=1}^l \alpha_i \mathbf{x}_i.$$

- 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

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

- Note: XX' may be close to singular, or singular according to machine precision.

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 λ defines the trade-off between norm and loss. This controls the complexity of the model (the process is called *regularization*).

- Taking the derivative and making it equal to zero:

$$\mathbf{X}'\mathbf{X}\mathbf{w} + \lambda\mathbf{w} = (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_n)\mathbf{w} = \mathbf{X}'\mathbf{y},$$

where \mathbf{I}_n is an identity matrix of $n \times n$ dimension,

- then,

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

- In terms of α :

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

- then

$$\alpha = \lambda^{-1}(\mathbf{y} - \mathbf{X}\mathbf{w}) = (\mathbf{X}\mathbf{X}' + \lambda\mathbf{I}_l)^{-1}\mathbf{y}.$$

Prediction function

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

Introduction

The Kernel
Trick

The Kernel
Approach to
Machine
Learning

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 $k_{ij} = k(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$
- The α 's are calculated as:

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

- The predicted function is approximated as:

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

Theorem

(Mercer's Theorem)

A function

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

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

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

into a feature map ϕ 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.

Some kernel functions

Introduction

The Kernel
Trick

The Kernel
Approach to
Machine
Learning

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Mathematical
characterisation

Visualizing kernels in
input space

Kernel
Algorithms

Kernels in
Complex
Structured
Data

Assume k_1 and k_2 kernels:

- $k(x, z) = p(k_1(x, z))$. p a polynomial with positive coefficients.
- $k(x, z) = \exp(k_1(x, z))$.
- $k(x, z) = \exp(-\|x - z\|^2 / (2\sigma^2))$. Gaussian kernel.
- $k(x, z) = k_1(x, z)k_2(x, z)$

Embeddings corresponding to kernels

Introduction

The Kernel
Trick

The Kernel
Approach to
Machine
Learning

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

**Mathematical
characterisation**

Visualizing kernels in
input space

Kernel
Algorithms

Kernels in
Complex
Structured
Data

- 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?

Introduction

The Kernel
Trick

The Kernel
Approach to
Machine
Learning

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Mathematical
characterisation

Visualizing kernels in
input space

Kernel
Algorithms

Kernels in
Complex
Structured
Data

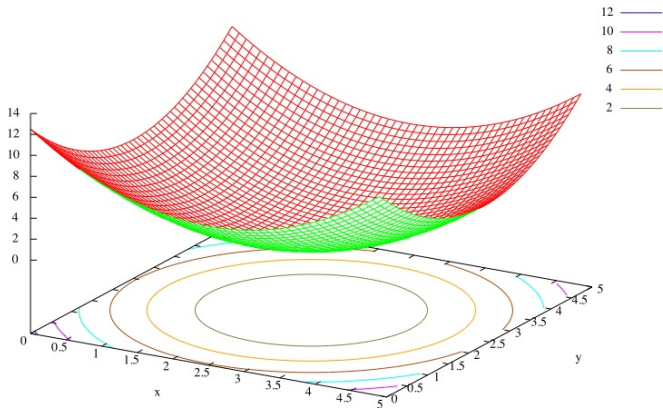
- Choose a point in input space p_0
- Calculate the distance from another point x to p_0 in the feature space:

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

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

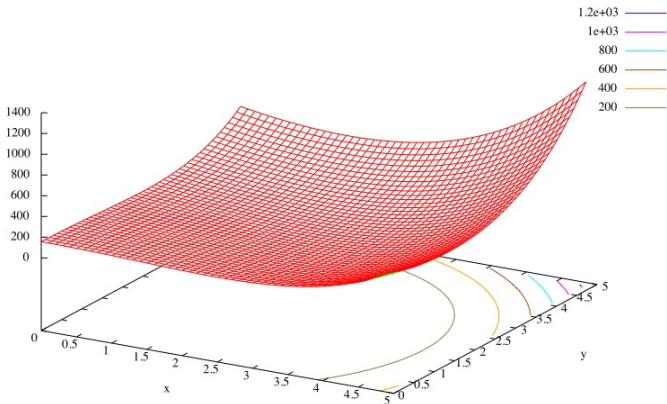
Identity kernel

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



Quadratic kernel (1)

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



Introduction

The Kernel
Trick

The Kernel
Approach to
Machine
Learning

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Mathematical
characterisation

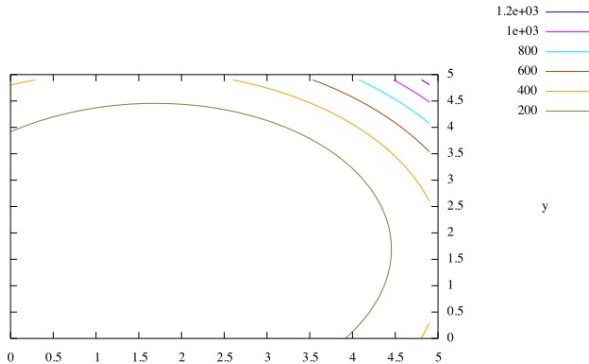
Visualizing kernels in
input space

Kernel
Algorithms

Kernels in
Complex
Structured
Data

Identity kernel (2)

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



Introduction

The Kernel
Trick

The Kernel
Approach to
Machine
Learning

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Mathematical
characterisation

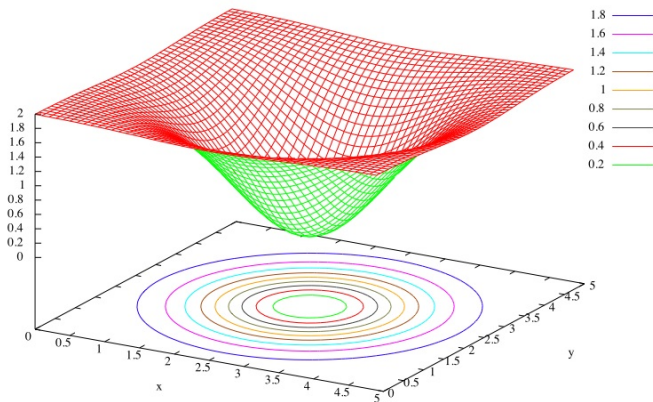
Visualizing kernels in
input space

Kernel
Algorithms

Kernels in
Complex
Structured
Data

Gaussian kernel

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



Introduction

The Kernel
Trick

The Kernel
Approach to
Machine
Learning

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Mathematical
characterisation

Visualizing kernels in
input space

Kernel
Algorithms

Kernels in
Complex
Structured
Data

Basic computations in feature space

- Means
- Distances
- Projections
- Covariance

Introduction

The Kernel
Trick

The Kernel
Approach to
Machine
Learning

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

Kernels in
Complex
Structured
Data

Classification and regression

Introduction

The Kernel
Trick

The Kernel
Approach to
Machine
Learning

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

**Kernel
Algorithms**

Kernels in
Complex
Structured
Data

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

Dimensionality reduction and clustering

Introduction

The Kernel
Trick

The Kernel
Approach to
Machine
Learning

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

**Kernel
Algorithms**

Kernels in
Complex
Structured
Data

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

Kernels in complex structured data

Introduction

The Kernel
Trick

The Kernel
Approach to
Machine
Learning

A Kernel
Pattern
Analysis
Algorithm

Kernel
Functions

Kernel
Algorithms

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

How to do symbolic regression?

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

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

Solution

- Define a kernel on strings

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

- 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

-

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

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

- $k(\text{ACTAG}, \text{CCTCG}) = ?$

- Is it a kernel?

Induced Feature Space

- What is the feature space induced by k ?

-

$$\begin{aligned}\phi : \Sigma^d &\rightarrow \mathbb{R}^{4d} \\ s_1 \dots s_d &\mapsto (x_1^1, \dots, x_4^1, x_1^2, \dots, x_4^2, \dots, x_1^d, \dots, x_4^d)\end{aligned}$$

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

References



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