Introduction to Kernel Methods (part I)

Liva Ralaivola

liva@ics.uci.edu

School of Information and Computer Science Institute for Genomics and Bioinformatics

Outline

- Learning to classify
- Vectors and inner products
- A simple linear classifier
- *Kernel trick* and linear classifiers in *feature spaces*
- Common kernels

Classifying things (1/2)

Classification is a real-world task

- does some patient have a serious disease?
- what kind of secondary structure does a sequence of amino acids correspond to?
- is some molecule toxic/not toxic for some organism?
- Automating this task
 - dealing with large number of items to classify
 - speed
 - cost

Classifying things (2/2)

Important notions in *learning to classify*

- limited number of *training* data (patients, sequences, molecules, etc.)
- learning algorithm (how to build the classifier?)
- generalization: the classifier should correctly classify test data
- Quick formalization
 - \mathcal{X} (e.g. $\mathbb{R}^d, d > 0$) is the space of data, called *input space*
 - \mathcal{Y} (e.g. toxic/not toxic, or $\{-1,+1\}$) is the target space
 - $f: \mathcal{X} \to \mathcal{Y}$ is the classifier



Vectors and inner product (1/3)



Vectors and inner product (2/3)



Vectors and inner product (3/3)



⟨u - v, e⟩ > 0: u - v and e point to the 'same direction'
⟨u - v, f⟩ = 0: u - v and f are orthogonal
⟨u - v, g⟩ < 0: u - v and g point to 'opposite directions'

A simple linear classifier



- Idea (see [Schölkopf and Smola, 2002] for details): classify points x according to which of the two class means c⁺ or c⁻ is closer:
 - for $x \in \mathcal{X}$, it is sufficient to take the sign of the inner product between w and x c

If $h(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} - \mathbf{c} \rangle$, we have the classifier $f(\mathbf{x}) = \operatorname{sign} (h(\mathbf{x}))$

• the (dotted) hyperplane (H), of normal vector \mathbf{w} , is the decision surface

A simple linear classifier



$$h(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} - \mathbf{c} \rangle = \langle \mathbf{w}, \mathbf{x} \rangle - \langle \mathbf{w}, \mathbf{c} \rangle$$
$$= \langle \mathbf{c}^+, \mathbf{x} \rangle - \langle \mathbf{c}^-, \mathbf{x} \rangle - \langle \mathbf{c}^+, \mathbf{c} \rangle + \langle \mathbf{c}^-, \mathbf{c} \rangle$$
$$= \sum_{i=1,...,m} \alpha_i \langle \mathbf{x}_i, \mathbf{x} \rangle + b, \quad \text{with } b \text{ a real constant}$$

A simple linear classifier



The kernel trick (1/4)

- Context: nonlinearly separable dataset $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$
- Idea to learn a nonlinear classifier
 - \blacksquare choose a (nonlinear) mapping ϕ

$$\begin{array}{rcccc} \phi : & \mathcal{X} & \to & \mathcal{H} \\ & \mathbf{x} & \mapsto & \phi(\mathbf{x}) \end{array}$$

where \mathcal{H} is an inner product space (inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$), called *feature space*

find a linear classifier (i.e. a separating hyperplane) in \mathcal{H} to classify $\{(\phi(\mathbf{x}_1), y_1), \dots, (\phi(\mathbf{x}_m), y_m)\}$

The kernel trick (2/4)



Taking the previous linear algorithm and implementing it in \mathcal{H} :

$$h(\mathbf{x}) = \sum_{i=1,\dots,m} \alpha_i \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle_{\mathcal{H}} + b$$

The kernel trick (3/4)

The kernel trick can be applied if there is a function k : X × X → ℝ such that: k(u, v) = ⟨φ(u), φ(v)⟩_H If so, all occurrences of ⟨φ(x_i), φ(x)⟩_H are replaced by k(x_i, x)

Keypoint: the 'focus' is sometimes only on k and not on ϕ

- Kernels must verify Mercer's property to be valid kernels
 - ensures that there exist a space \mathcal{H} and a mapping $\phi : \mathcal{X} \to \mathcal{H}$ such that $k(\mathbf{u}, \mathbf{v}) = \langle \phi(\mathbf{u}), \phi(\mathbf{v}) \rangle_{\mathcal{H}}$

however non valid kernels have been used with success

- and, research is in progress on using non semi-definite kernels
- \blacksquare k might be viewed as a similarity measure

The kernel trick (4/4)



consider a nonlinear classification problem on $\mathcal{X} \times \mathcal{Y}$

• choose a linear classification algorithm (expr. in terms $\langle \cdot, \cdot \rangle$)

replace all occurrences of $\langle \cdot, \cdot \rangle$ by a kernel $k(\cdot, \cdot)$

Obtained classifier:
$$f(\mathbf{x}) = \operatorname{sign} \left(\sum_{i=1,...,m} \alpha_i k(\mathbf{x}_i, \mathbf{x}) + b \right)$$

Common kernels (1/2)



Common kernels (2/2)

- Let $k = \langle \mathbf{u}, \mathbf{v} \rangle_{\mathbb{R}^2}^2$ (polynomial kernel with c = 0 and d = 2) defined on $\mathbb{R}^2 \times \mathbb{R}^2$
- Consider the mapping:

$$\phi: \quad \mathbb{R}^2 \quad \to \quad \mathbb{R}^3$$
$$\mathbf{x} = [x_1, x_2]^\top \quad \mapsto \quad \phi(\mathbf{x}) = \left[x_1^2, \sqrt{2}x_1 x_2, x_2^2\right]^\top$$

• We have, for $\mathbf{u}, \mathbf{v} \in \mathbb{R}^2$:

$$\begin{split} \langle \phi(\mathbf{u}), \phi(\mathbf{v}) \rangle_{\mathbb{R}^3} &= \langle [u_1^2, \sqrt{2}u_1u_2, u_2^2]^\top, [v_1^2, \sqrt{2}v_1v_2, v_2^2]^\top \rangle \\ &= (u_1v_1 + u_2v_2)^2 \\ &= \langle \mathbf{u}, \mathbf{v} \rangle_{\mathbb{R}^2}^2 \\ &= k(\mathbf{u}, \mathbf{v}) \end{split}$$

What to take home ?

Classification

appears in many real world problems

can be done automatically

generalization is important

Kernels

are wonderful !

field of active, thrilling research

classification of structured objects might be envisioned:

sequences: DNA strings, amino-acid strings, texts

 \Rightarrow Kernels and appl. to sequences

graphs: structure of a molecule, disulfide bonds, GO

 \Rightarrow Kernels and appl. to graphs and molecular structures

fusion of heterogeneous information

 \Rightarrow Combining Classifiers and Combining Kernels

Miscellaneous information

- Contact information
 - Liva Ralaivola
 - Office CS/E 307
 - Webpage (not really up to date): http://www.ics.uci.edu/~liva/
- Useful places/resources/links
 - Kernel Club: every Friday starting Oct. 1st, 2004, 11:00am, CS 432
 - Kernel Machines: http://www.kernel-machines.org/
 - Learning with Kernels: http://www.learning-with-kernels.org/
 - SVM applet: http://svm.dcs.rhbnc.ac.uk/pagesnew/GPat.shtml

References

References

[Schölkopf and Smola, 2002] Schölkopf, B. and Smola, A. J. (2002). Learning with Kernels, Support Vector Machines, Regularization, Optimization and Beyond. MIT University Press.

Other references on my web page (soon)