

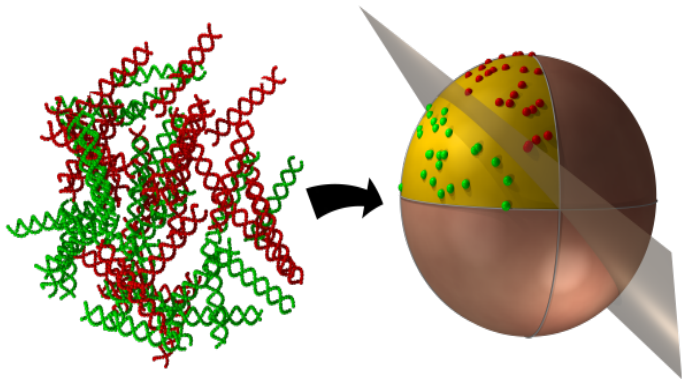
Ch7. Kernel Methods

ST4240, 2014/2015

Version 0.1

Alexandre Thiéry

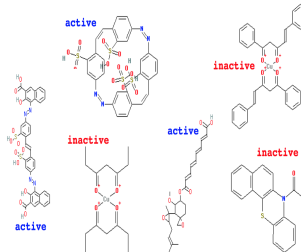
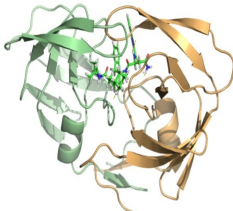
Department of Statistics and Applied Probability



MASKATLLLAFTLLFATCIARHQQRQQQNQCQLQNIEA...

MARSSLFTFLCLAVFINGCLSQIEQQSPWEFQGSEVW...

MALHTVLIIMLSLLPMLEAQNPEHANITIGEPITNETLGWL...



0	4	1	9	2	1	3	1	4	3
5	3	6	1	7	2	8	6	9	4
0	9	1	1	2	4	3	2	7	3
8	6	9	0	5	6	0	7	6	1
8	7	9	3	9	8	5	9	3	3
0	7	4	9	8	0	9	4	1	4
4	6	0	4	5	6	1	0	0	1
7	1	6	3	0	2	1	1	7	9
0	2	6	7	8	3	9	0	4	6
7	4	6	8	0	7	8	3	1	5

Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based essentially on the messages that reach our eyes. For a long time, the visual image was considered as a

image was considered as a visual centers in the brain. As a movie scene, the image is a discovery. The knowledge of perception is more complex. Following the path to the various parts of the cortex, Hubel and Wiesel have demonstrated that the message about the image falling on the retina undergoes a

analysis in a system of nerve cells stored in columns. In this system each cell has its specific function and is responsible for a specific detail in the pattern of the retinal image.

**sensory, brain,
visual, perception,
retinal, cerebral cortex,
eye, cell, optical
nerve, image
Hubel, Wiesel**

China is forecasting a trade surplus of \$90bn (£51bn) to \$100bn this year, a threefold increase on 2004's \$32bn. The Commerce Ministry said the surplus would be created by a predicted 30% increase in exports to \$750bn, compared with \$575bn in 2004.

\$660bn. The increase will annoy the US, which China's deliberate policy to agree to a yuan is a government also needs to demand some country. China's yuan against the dollar and permitted it to trade within a narrow band, but the US wants the yuan to be allowed to move freely. However, Beijing has made it clear that it will take its time and tread carefully before allowing the yuan to rise further in value.

**China, trade,
surplus, commerce,
exports, imports, US,
yuan, bank, domestic,
foreign, increase,
trade, value**

Outline

1 General Overview

2 Kernels

3 From kernels to functions

Regularization and Optimization

- Training examples: $(x_i, y_i)_{i=1}^N$ with $x_i \in \mathcal{X}$ and $y_i \in \mathcal{Y}$
- Class of functions \mathcal{C} , regularisation functional $\Omega(\cdot)$

$$\operatorname{argmin}_{f \in \mathcal{C}} \left\{ f \mapsto \sum_{i=1}^N \mathbf{Loss}(f(x_i), y_i) + \lambda \cdot \Omega(f) \right\}$$

- **[Exercise]** : OLS, Ridge regression, LASSO, Logistic regression, SVM, Boosting?

$$\operatorname{argmin}_{f \in \mathcal{C}} \left\{ f \mapsto \sum_{i=1}^N \mathbf{Loss}(f(x_i), y_i) + \lambda \cdot \Omega(f) \right\}$$

- Linear models: $f(x) = \langle \beta, x \rangle$ with $\beta \in \mathbb{R}^p$.
- What if the class of function \mathcal{C} is very large, or infinite dimensional?
- **Example:** regression with $\mathcal{C} = (\text{Smooth functions})$ and

$$\Omega(f) = \int \|f'\|^2(u) du.$$

Unstructured Data and feature extraction

- Digit classification: what features?
- Spam v.s. No Spam: what features?
- Molecule Classification: what features?

Objective

- For a general set \mathcal{X} , define a sensible class of functions \mathcal{C}
- Define the regularisation functional Ω such that the optimisation problem

$$\operatorname{argmin}_{f \in \mathcal{C}} \left\{ f \mapsto \sum_{i=1}^N \mathbf{Loss}(f(x_i), y_i) + \lambda \cdot \Omega(f) \right\}$$

can be efficiently solved, even if \mathcal{C} is huge.

- Understand why a large class of function \mathcal{C} does not necessarily lead to bad performances (overfitting?)

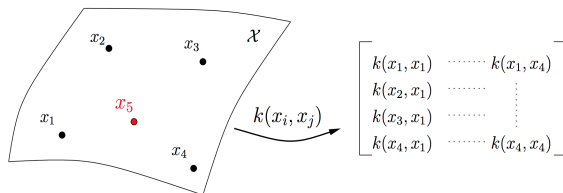
Outline

1 General Overview

2 Kernels

3 From kernels to functions

Measure of similarity



- Unstructured data $\{x_i\}_{i=1}^N$
- Use a **symmetric kernel** $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ to measure similarity

$$k(x_i, x_j) \approx (\text{similarity between } x_i \text{ and } x_j)$$

- Consider the **Gram matrix** \mathbf{K} defined as

$$[K]_{i,j} = k(x_i, x_j)$$

- \mathbf{K} contains all the pairwise measures of similarity

- The kernel **kernel** $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is positive semidefinite if for any $\{x_i\}_{i=1}^N$ the **Gram matrix** \mathbf{K} defined as

$$[K]_{i,j} = k(x_i, x_j)$$

is **positive semi-definite**.

- **[Exercise]** What about $\mathcal{X} = \mathbb{R}^d$ and $k(x_i, x_j) = \langle x_i, x_j \rangle$?
- **[Exercise]** For an arbitrary set \mathcal{X} and a **feature map** $\varphi : \mathcal{X} \rightarrow \mathbb{R}^d$, what about

$$k(x_i, x_j) = \langle \varphi(x_i), \varphi(x_j) \rangle?$$

Operation on kernels

- Suppose that k_1 and k_2 are two positive semidefinite kernels on \mathcal{X} .
- **[Exercise]** : is $k_+(x_i, x_j) = k_1(x_i, x_j) + k_2(x_i, x_j)$ positive semidefinite?
- **[Exercise]** : is $k_\times = k_1(x_i, x_j) \times k_2(x_i, x_j)$ positive semidefinite?
- **[Exercise]** : is $k_{\exp}(x_i, x_j) = \exp(k(x_i, x_j))$ positive semidefinite?
- **[Exercise]** : is $k_{poly}(x_i, x_j) = (1 + k(x_i, x_j))^p$ positive semidefinite?
- **Radial Basis Function kernel**: the kernel on \mathbb{R}^d defined as

$$k_{RBF}(x_i, x_j) = \exp \left\{ -\frac{\|x_i - x_j\|^2}{\ell^2} \right\}$$

is positive definite.

String kernel

IPTSALVKETLALLSTHRTLLIANETLRIPVPVHKNHQLCTEEIFQGIGTLESQTVQGGTV
ERLFKNLSLIKKYIDGQKKKCGEERRRVNQFLDYLQEFLGVMNTEWI

PHRRDLCSRSIWLARKIRSDLTALTESYVKHQGLWSELTEAERLQENLQAYRTFHVLLA
RLLEDQQVHFTPTGDFHQAHTLLQVAAFAYQIEELMILLEYKIPRNEADGMLFEKK
LWGLKVLQELSQWTVRSIHDLRFISSHQTGIP

- One can define a kernel on strings of character. For two strings x and y define

$$\mathbf{k}(x, y) = \sum_{s \in \mathcal{S}} \omega_s \varphi_s(x) \varphi_s(y)$$

where:

- \mathcal{S} a set of possible substring
- ω_s is a positive weight
- $\varphi_s(x)$ is the number of occurrences of s in x .
- **[Exercise]**: why is it positive semi-definite?

Outline

1 General Overview

2 Kernels

3 From kernels to functions

Reproducing Kernel Hilbert Space

$$\operatorname{argmin}_{f \in \mathcal{C}} \left\{ f \mapsto \sum_{i=1}^N \mathbf{Loss}(f(x_i), y_i) + \lambda \cdot \Omega(f) \right\}$$

- Recall that one were interested in defining a class \mathcal{C} of functions on \mathcal{X} and a regularisation functional Ω
- Suppose that $k(\cdot, \cdot)$ is a positive semi-definite kernel on \mathcal{X} .
- For any $x_0 \in \mathcal{X}$ one can define a function $F_{x_0}(\cdot) : \mathcal{X} \rightarrow \mathbb{R}$ by

$$\mathbf{F}_{\mathbf{x}_0}(\mathbf{x}) = \mathbf{k}(\mathbf{x}_0, \mathbf{x})$$

Reproducing Kernel Hilbert Space

- The RKHS \mathcal{H} of function is the class \mathcal{C} of functions on \mathcal{X} that can be expressed as

$$F(\cdot) = \sum_i \alpha_i F_{x_i}(\cdot)$$

for some elements $\{x_i\}_{i \in I}$ of \mathcal{X}

- Note that \mathcal{H} vector space of functions
- One can define the norm

$$\|F\|_{RKHS}^2 = \sum_{i,j} \alpha_i \alpha_j k(x_i, x_j)$$

Representer theorem

- Class of functions: $\mathcal{C} \equiv \mathcal{H}$
- Regularization: $\Omega(f) \equiv \Psi(\|f\|_{RKHS})$ for strictly increasing Ψ
- In other words

$$\operatorname{argmin}_{f \in \mathcal{H}} \left\{ f \mapsto \sum_{i=1}^N \mathbf{Loss}(f(x_i), y_i) + \lambda \cdot \Psi(\|f\|_{\mathcal{H}}) \right\} \quad (1)$$

Theorem (Representer Theorem)

A solution f_ to the optimisation problem (1) can always be expressed as*

$$f_*(\cdot) = \sum_{i=1}^N \alpha_i F_{x_i}(\cdot)$$

where $\{(x_i, y_i)\}_{i=1}^N$ are the training examples and $\{\alpha_i\}_{i=1}^N$ are some coefficients.

Representer theorem: consequence

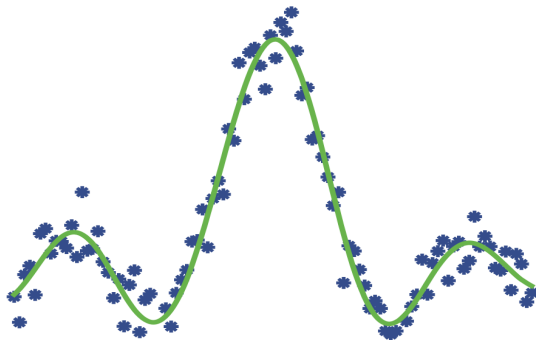
- Original optimisation problem

$$\operatorname{argmin}_{f \in \mathcal{H}} \left\{ f \mapsto \sum_{i=1}^N \mathbf{Loss}(f(x_i), y_i) + \lambda \cdot \Psi(\|f\|_{\mathcal{H}}) \right\} \quad (2)$$

- Thanks to the representer theorem, optimisation over $f \in \mathcal{H}$ reduces to optimisation over $\alpha \in \mathbb{R}^N$

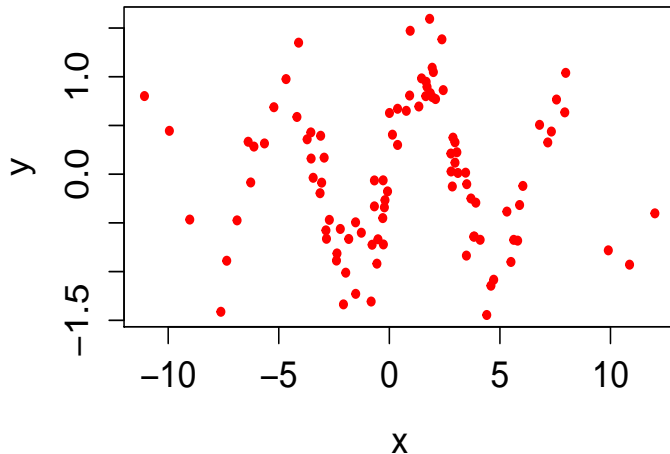
$$\operatorname{argmin}_{\alpha \in \mathbb{R}^N} \left\{ \alpha \mapsto \sum_{i=1}^N \mathbf{Loss}(f(x_i), [K\alpha]_i) + \lambda \cdot \langle \alpha, K\alpha \rangle \right\} \quad (3)$$

Kernel Ridge Regression

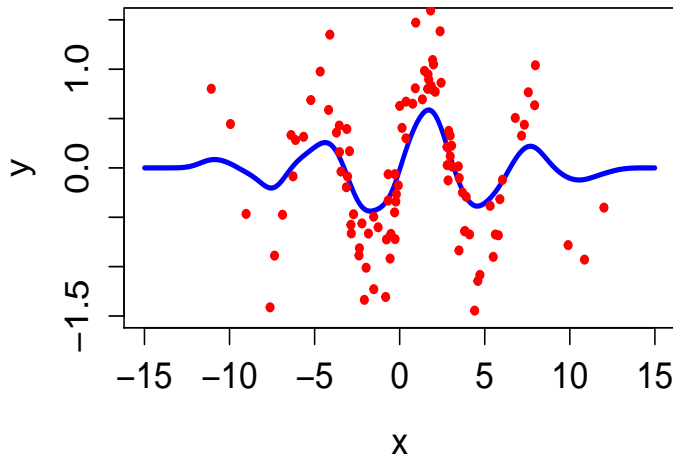


- **[Exercise]** What about ridge regression with kernel $k(\cdot, \cdot)$?

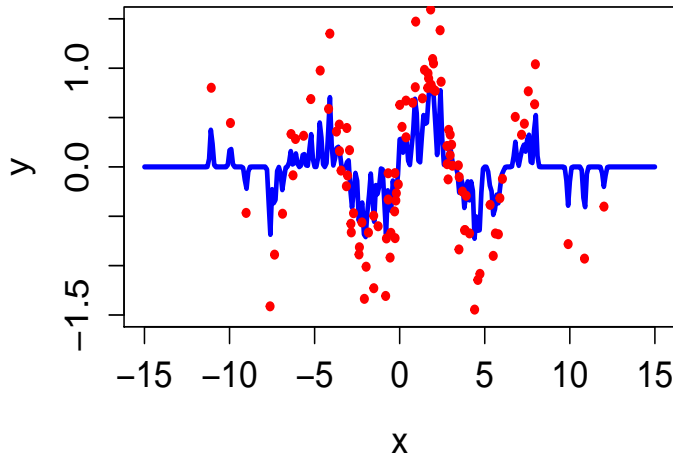
Data to be fitted



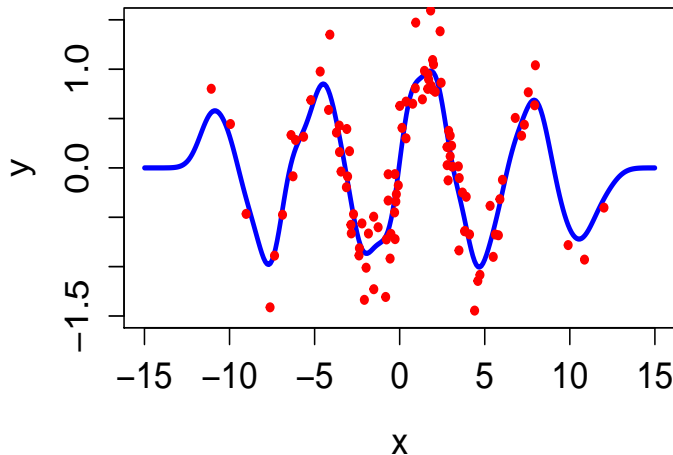
Kernel ridge $L=1$, $\lambda=10$



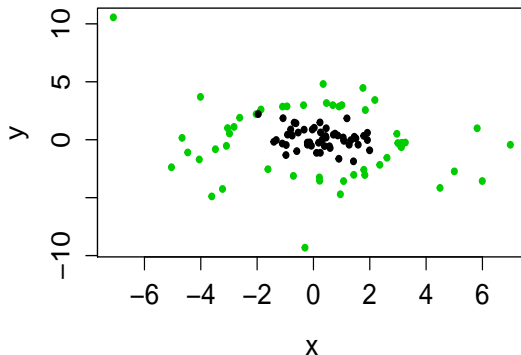
Kernel ridge $L = 0.1$, $\lambda = 1$



Kernel ridge $L=1$, $\lambda=0.5$

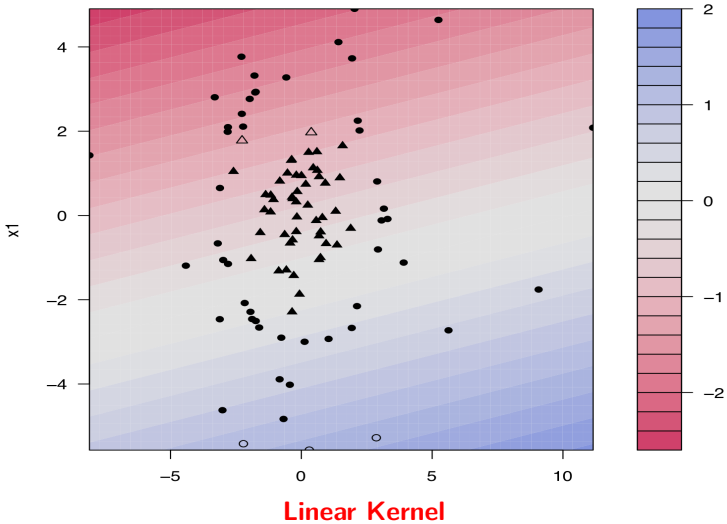


non linearly separable dataset

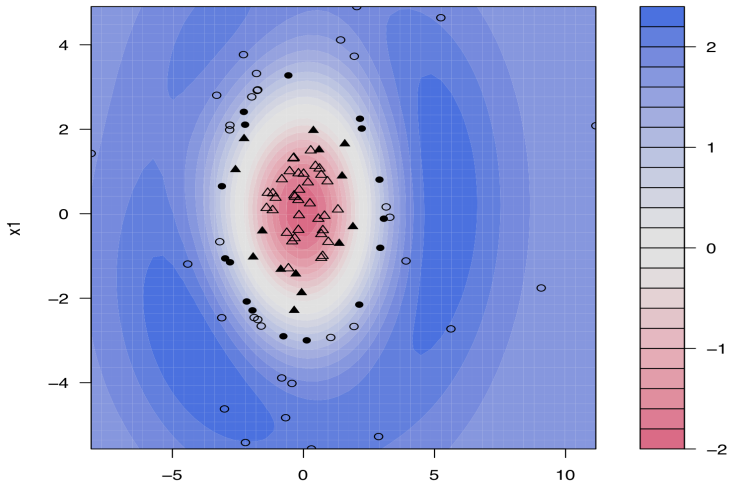


- **[Exercise]** What SVM with a kernel $k(\cdot, \cdot)$?

Kernelized SVM

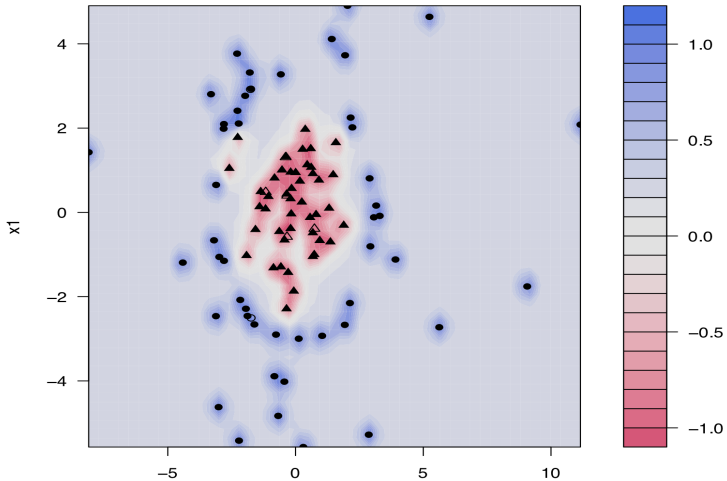


Kernelized SVM



RBF kernel with $\sigma = 1/\ell = 0.5$

Kernelized SVM



RBF kernel with $\sigma = 1/\ell = 50$