

Apprentissage, Noyaux et parcimonie

Janvier, 2009

Stat Im'09

Laboratoire d'Informatique, de Traitement de l'Information et des Systèmes

Stéphane Canu

stephane.canu@litislab.eu



Roadmap

1 Kernels and the learning problem

- Two learning problems
- Kernelizing the linear regression
- Kernel machines: a definition

2 Tools: the functional framework

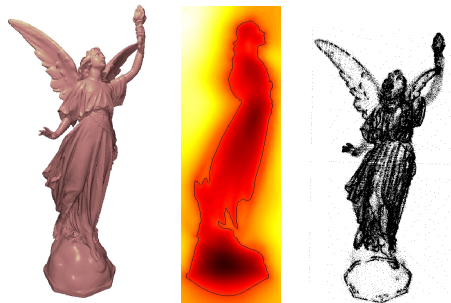
- In the beginning was the kernel
- Kernel and hypothesis set

3 Kernel machines

- Non sparse kernel machines
- sparse kernel machines: SVM
- practical SVM

4 Conclusion

Implicit Surface Modelling

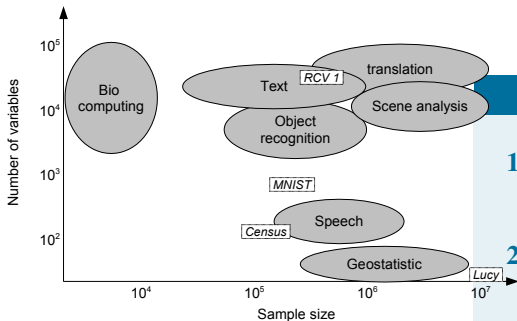


Example (*Lucy*²)

- ▶ 14 million points with normals.
- ▶ 364,982 compactly supported basis function centres.
- ▶ kernel regression with *thin plate splines* kernels

²<http://www.kyb.mpg.de/~walder>

Learning challenges: the size effect



3 key issues

1. learn any problem:
 - ▶ **functional universality**
2. from data:
 - ▶ **statistical consistency**
3. with large data sets:
 - ▶ **computational efficiency**

kernel machines adress these three issues
(up to a certain point regarding efficiency)

the example of the linear least mean square

the linear model

$$y_i = \sum_{j=1}^d \beta_j x_{ij} + \varepsilon_i \quad , \quad i = 1, n$$

n observations and d variables; $d < n$

$$\min_{\beta} = \sum_{i=1}^n \left(\sum_{j=1}^d x_{ij} \beta_j - y_i \right)^2 = \|X\beta - Y\|^2$$

Solution: $\hat{\beta} = (X^T X)^{-1} X^T Y$

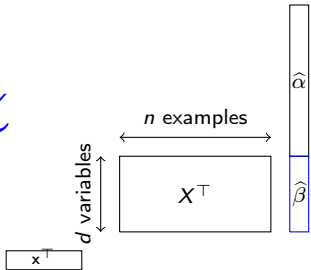
$$f(\mathbf{x}) = \mathbf{x}^T \underbrace{(X^T X)^{-1} X^T Y}_{\hat{\beta}}$$

What is the influence of each example (X rows)?

The influence of the examples

for a new input \mathbf{x}

$$\begin{aligned}
 f(\mathbf{x}) &= \mathbf{x}^\top (X^\top X)(X^\top X)^{-1} \underbrace{(X^\top X)^{-1} X^\top Y}_{\hat{\boldsymbol{\beta}}} \\
 &= \mathbf{x}^\top X^\top \underbrace{X(X^\top X)^{-1} (X^\top X)^{-1} X^\top Y}_{\hat{\boldsymbol{\alpha}}}
 \end{aligned}$$



$$f(\mathbf{x}) = \sum_{j=1}^d \hat{\beta}_j x_j$$

from variables to examples

$$\underbrace{\hat{\boldsymbol{\alpha}} = X(X^\top X)^{-1} \hat{\boldsymbol{\beta}}}_{n \text{ examples}}$$

and

$$\underbrace{\hat{\boldsymbol{\beta}} = X^\top \hat{\boldsymbol{\alpha}}}_{d \text{ variables}}$$

What if $d \geq n$?

Non linear case: dictionary vs. kernel

in the non linear case: use a **dictionary** of functions

$$\phi_j(\mathbf{x}), j = 1, p \quad \text{with possibly } p = \infty$$

for instance polynomials, wavelets... (assume orthogonality)

$$f(\mathbf{x}) = \sum_{j=1}^p \hat{\beta}_j \phi_j(\mathbf{x}) \quad \text{with} \quad \hat{\beta}_j = \sum_{i=1}^n y_i \phi_j(\mathbf{x}_i)$$

using linearity

$$f(\mathbf{x}) = \sum_{i=1}^n y_i \underbrace{\sum_{j=1}^p \phi_j(\mathbf{x}_i) \phi_j(\mathbf{x})}_{k(\mathbf{x}_i, \mathbf{x})}$$

$p \geq n$ so what since $k(\mathbf{x}_i, \mathbf{x}) = \sum_{j=1}^p \phi_j(\mathbf{x}_i) \phi_j(\mathbf{x})$

Non linear case: dictionary vs. kernel

in the non linear case: use a **dictionary** of functions

$$\phi_j(\mathbf{x}), j = 1, p \quad \text{with possibly} \quad p = \infty$$

for instance polynomials, wavelets... (assume orthogonality)

$$f(\mathbf{x}) = \sum_{j=1}^p \hat{\beta}_j \phi_j(\mathbf{x}) \quad \text{with} \quad \hat{\beta}_j = \sum_{i=1}^n y_i \phi_j(\mathbf{x}_i)$$

using linearity

$$f(\mathbf{x}) = \sum_{i=1}^n y_i \underbrace{\sum_{j=1}^p \phi_j(\mathbf{x}_i) \phi_j(\mathbf{x})}_{k(\mathbf{x}_i, \mathbf{x})}$$

$$p \geq n \text{ so what since } k(\mathbf{x}_i, \mathbf{x}) = \sum_{j=1}^p \phi_j(\mathbf{x}_i) \phi_j(\mathbf{x})$$

closed form kernel: the quadratic kernel

The quadratic dictionary in \mathbb{R}^d :

$$\begin{aligned}\Phi : \mathbb{R}^d &\rightarrow \mathbb{R}^{p=1+d+\frac{d(d+1)}{2}} \\ \mathbf{s} &\mapsto \Phi = (1, s_1, s_2, \dots, s_d, s_1^2, s_2^2, \dots, s_d^2, \dots, s_i s_j, \dots)\end{aligned}$$

in this case

$$\Phi(\mathbf{s})^\top \Phi(\mathbf{t}) = 1 + s_1 t_1 + s_2 t_2 + \dots + s_d t_d + s_1^2 t_1^2 + \dots + s_d^2 t_d^2 + \dots + s_i s_j t_i t_j + \dots$$

The quadratic kernel: $\mathbf{s}, \mathbf{t} \in \mathbb{R}^d$, $k(\mathbf{s}, \mathbf{t}) = (\mathbf{s}^\top \mathbf{t} + 1)^2 = 1 + 2\mathbf{s}^\top \mathbf{t} + (\mathbf{s}^\top \mathbf{t})^2$

computes the dot product of the reweighted dictionary:

$$\begin{aligned}\Phi : \mathbb{R}^d &\rightarrow \mathbb{R}^{p=1+d+\frac{d(d+1)}{2}} \\ \mathbf{s} &\mapsto \Phi = (1, \sqrt{2}s_1, \sqrt{2}s_2, \dots, \sqrt{2}s_d, s_1^2, s_2^2, \dots, s_d^2, \dots, \sqrt{2}s_i s_j, \dots)\end{aligned}$$

$p = 1 + d + \frac{d(d+1)}{2}$ multiplications vs. $d + 1$
use kernel to save computation

closed form kernel: the quadratic kernel

The quadratic dictionary in \mathbb{R}^d :

$$\begin{aligned} \Phi : \mathbb{R}^d &\rightarrow \mathbb{R}^{p=1+d+\frac{d(d+1)}{2}} \\ \mathbf{s} &\mapsto \Phi = (1, s_1, s_2, \dots, s_d, s_1^2, s_2^2, \dots, s_d^2, \dots, s_i s_j, \dots) \end{aligned}$$

in this case

$$\Phi(\mathbf{s})^\top \Phi(\mathbf{t}) = 1 + s_1 t_1 + s_2 t_2 + \dots + s_d t_d + s_1^2 t_1^2 + \dots + s_d^2 t_d^2 + \dots + s_i s_j t_i t_j + \dots$$

The quadratic kernel: $\mathbf{s}, \mathbf{t} \in \mathbb{R}^d, \quad k(\mathbf{s}, \mathbf{t}) = (\mathbf{s}^\top \mathbf{t} + 1)^2$
 $= 1 + 2\mathbf{s}^\top \mathbf{t} + (\mathbf{s}^\top \mathbf{t})^2$

computes the dot product of the reweighted dictionary:

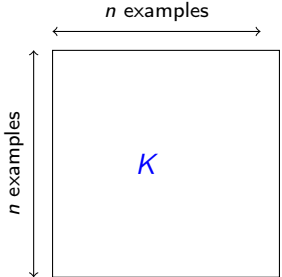
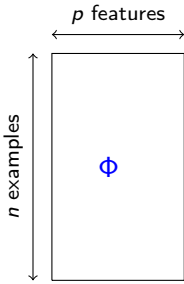
$$\begin{aligned} \Phi : \mathbb{R}^d &\rightarrow \mathbb{R}^{p=1+d+\frac{d(d+1)}{2}} \\ \mathbf{s} &\mapsto \Phi = (1, \sqrt{2}s_1, \sqrt{2}s_2, \dots, \sqrt{2}s_d, s_1^2, s_2^2, \dots, s_d^2, \dots, \sqrt{2}s_i s_j, \dots) \end{aligned}$$

$p = 1 + d + \frac{d(d+1)}{2}$ multiplications vs. $d + 1$
 use kernel to save computation

kernel: features through pairwise comparizons

\mathbf{x}
 e.g. a text

$\phi(\mathbf{x})$
 e.g. BOW



$$k(\mathbf{x}_i, \mathbf{x}_j) = \sum_{j=1}^p \phi_j(\mathbf{x}_i)\phi_j(\mathbf{x}_j)$$

K The matrix of *pairwise comparizons* ($\mathcal{O}(n^2)$)

Kernel machines

use a kernel as a dictionary

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i)$$

- ▶ α_i influence of example i depends on y_i
- ▶ $k(\mathbf{x}, \mathbf{x}_i)$ the kernel do NOT depend on y_i

Definition (Kernel)

a function k from $\mathcal{X} \times \mathcal{X}$ onto \mathbb{R} .

semi-parametric version: given the family $q_j(\mathbf{x}), j = 1, p$

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i) + \sum_{j=1}^p \beta_j q_j(\mathbf{x})$$

Kernel machines

use a kernel as a dictionary

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i)$$

- ▶ α_i influence of example i depends on y_i
- ▶ $k(\mathbf{x}, \mathbf{x}_i)$ the kernel do NOT depend on y_i

Definition (Kernel)

a function k from $\mathcal{X} \times \mathcal{X}$ onto \mathbb{R} .

semi-parametric version: given the family $q_j(\mathbf{x}), j = 1, p$

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i) + \sum_{j=1}^p \beta_j q_j(\mathbf{x})$$

Kernel machines

Definition (Kernel machines)

$$\mathcal{A}((\mathbf{x}_i, y_i)_{i=1, n})(\mathbf{x}) = \psi\left(\sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i) + \sum_{j=1}^p \beta_j q_j(\mathbf{x})\right)$$

α et β : parameters to be estimated.

Exemples

$$\mathcal{A}(x) = \sum_{i=1}^n \alpha_i (x - x_i)_+^3 + \beta_0 + \beta_1 x \quad \text{splines}$$

$$\mathcal{A}(\mathbf{x}) = \text{sign}\left(\sum_{i \in I} \alpha_i \exp^{-\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{b}} + \beta_0\right) \quad \text{SVM}$$

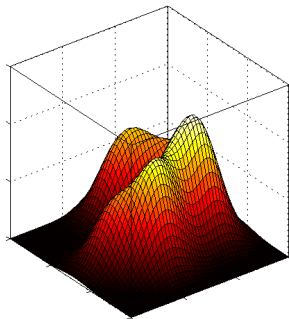
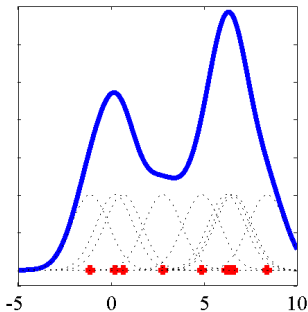
$$\mathbb{P}(y|\mathbf{x}) = \frac{1}{Z} \exp\left(\sum_{i \in I} \alpha_i \mathbb{1}_{\{y=y_i\}}(\mathbf{x}^\top \mathbf{x}_i + b)^2\right) \quad \text{exponential family}$$

example of kernel machine: the parzen estimate (1960)

assume the kernel is normalized: $\forall \mathbf{s} \in \mathcal{X}, \int_{\mathcal{X}} k(\mathbf{x}, \mathbf{s}) d\mathbf{x} = 1$

for a given data set $(\mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_n)$ the parzen window estimate is

$$\hat{\mathbb{P}}(x) = \frac{1}{n} \sum_{i=1}^n k(\mathbf{x}, \mathbf{x}_i)$$



- derive:
- potential based classification rule
 - the Nadaraya-Watson regression estimator

Different kernel machines

- ▶ different kernels
 - kernel machines are defined for (almost) ANY kernel
 - kernel machines are defined for (almost) ANY input
- ▶ kernel machines are defined through the cost functions
 - ▶ task dependent criterion:
 - ▶ classif (SVM, K Log. reg), regression (SVR, splines)
 - ▶ ranking, clustering (OCSVM), semi supervised (Trans. SVM)
 - ▶ dim. reduction (KPCA, KPLS), sources separation (KICA)...
 - ▶ penalty term
 - ▶ sparse / non sparse : $l_0 = \{\alpha_i = 0\}$

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i) + \sum_{j=1}^p \beta_j q_j(\mathbf{x})$$

- ▶ different implementations (algorithms)
 - introducing sparsity

linear algorithm → kernelization → sparsity

Roadmap

1 Kernels and the learning problem

- Two learning problems
- Kernelizing the linear regression
- Kernel machines: a definition

2 Tools: the functional framework

- In the beginning was the kernel
- Kernel and hypothesis set

3 Kernel machines

- Non sparse kernel machines
- sparse kernel machines: SVM
- practical SVM

4 Conclusion

In the beginning was the kernel...

Definition (Kernel)

a function of two variable k from $\mathcal{X} \times \mathcal{X}$ to \mathbb{R}

Definition (Positive kernel)

A kernel $k(s, t)$ on \mathcal{X} is said to be positive

- ▶ if it is symmetric: $k(s, t) = k(t, s)$
- ▶ and if for any finite positive integer n :

$$\forall \{\alpha_i\}_{i=1,n} \in \mathbb{R}, \forall \{\mathbf{x}_i\}_{i=1,n} \in \mathcal{X}, \quad \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 0$$

it is strictly positive if for $\alpha_i \neq 0$

$$\sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) > 0$$

Examples of positive kernels

the linear kernel: $\mathbf{s}, \mathbf{t} \in \mathbb{R}^d$, $k(\mathbf{s}, \mathbf{t}) = \mathbf{s}^\top \mathbf{t}$

symmetric: $\mathbf{s}^\top \mathbf{t} = \mathbf{t}^\top \mathbf{s}$

$$\begin{aligned} \text{positive: } \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) &= \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \mathbf{x}_i^\top \mathbf{x}_j \\ &= \left(\sum_{i=1}^n \alpha_i \mathbf{x}_i \right)^\top \left(\sum_{j=1}^n \alpha_j \mathbf{x}_j \right) = \left\| \sum_{i=1}^n \alpha_i \mathbf{x}_i \right\|^2 \end{aligned}$$

the product kernel: $k(\mathbf{s}, \mathbf{t}) = g(\mathbf{s})g(\mathbf{t})$ for some $g: \mathbb{R}^d \rightarrow \mathbb{R}$,

symmetric by construction

$$\begin{aligned} \text{positive: } \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) &= \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j g(\mathbf{x}_i) g(\mathbf{x}_j) \\ &= \left(\sum_{i=1}^n \alpha_i g(\mathbf{x}_i) \right) \left(\sum_{j=1}^n \alpha_j g(\mathbf{x}_j) \right) = \left(\sum_{i=1}^n \alpha_i g(\mathbf{x}_i) \right)^2 \end{aligned}$$

k is positive \Leftrightarrow (its square root exists) $\Leftrightarrow k(\mathbf{s}, \mathbf{t}) = \langle \phi_{\mathbf{s}}, \phi_{\mathbf{t}} \rangle$

positive definite Kernel (PDK) algebra (closure)

if $k_1(\mathbf{s}, \mathbf{t})$ and $k_2(\mathbf{s}, \mathbf{t})$ are two positive kernels

- ▶ DPK are a convex cone: $\forall a_1 \in \mathbb{R}^+ \quad a_1 k_1(\mathbf{s}, \mathbf{t}) + k_2(\mathbf{s}, \mathbf{t})$
- ▶ for any measurable function ψ from \mathcal{X} to \mathbb{R} $k(\mathbf{s}, \mathbf{t}) = \psi(\mathbf{s})\psi(\mathbf{t})$
- ▶ product kernel $k_1(\mathbf{s}, \mathbf{t})k_2(\mathbf{s}, \mathbf{t})$

proofs

▶ by linearity:

$$\sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j (a_1 k_1(i, j) k_2(i, j)) = a_1 \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k_1(i, j) + \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k_2(i, j)$$

▶ by linearity: $\sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j (\psi(\mathbf{x}_i) \psi(\mathbf{x}_j)) = \left(\sum_{i=1}^n \alpha_i \psi(\mathbf{x}_i) \right) \left(\sum_{j=1}^n \alpha_j \psi(\mathbf{x}_j) \right)$

▶ assuming $\exists \psi_\ell$ s.t. $k_1(\mathbf{s}, \mathbf{t}) = \sum_\ell \psi_\ell(\mathbf{s}) \psi_\ell(\mathbf{t})$

$$\begin{aligned} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k_1(\mathbf{x}_i, \mathbf{x}_j) k_2(\mathbf{x}_i, \mathbf{x}_j) &= \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \left(\sum_\ell \psi_\ell(\mathbf{x}_i) \psi_\ell(\mathbf{x}_j) k_2(\mathbf{x}_i, \mathbf{x}_j) \right) \\ &= \sum_\ell \sum_{i=1}^n \sum_{j=1}^n (\alpha_i \psi_\ell(\mathbf{x}_i)) (\alpha_j \psi_\ell(\mathbf{x}_j)) k_2(\mathbf{x}_i, \mathbf{x}_j) \end{aligned}$$

Kernel engineering: building PDK

- ▶ for any polynomial with positive coef. ϕ from \mathbb{R} to \mathbb{R}
 $\phi(k(\mathbf{s}, \mathbf{t}))$
- ▶ if Ψ is a function from \mathbb{R}^d to \mathbb{R}^d
 $k(\Psi(\mathbf{s}), \Psi(\mathbf{t}))$
- ▶ if φ from \mathbb{R}^d to \mathbb{R}^+ , is minimum in 0
 $k(\mathbf{s}, \mathbf{t}) = \varphi(\mathbf{s} + \mathbf{t}) - \varphi(\mathbf{s} - \mathbf{t})$
- ▶ convolution of two positive kernels is a positive kernel
 $K_1 \star K_2$

the Gaussian kernel is a PDK

$$\begin{aligned} \exp(-\|\mathbf{s} - \mathbf{t}\|^2) &= \exp(-\|\mathbf{s}\|^2 - \|\mathbf{t}\|^2 - 2\mathbf{s}^\top \mathbf{t}) \\ &= \exp(-\|\mathbf{s}\|^2) \exp(-\|\mathbf{t}\|^2) \exp(2\mathbf{s}^\top \mathbf{t}) \end{aligned}$$

- ▶ $\mathbf{s}^\top \mathbf{t}$ is a PDK and function \exp as the limit of positive series expansion, so $\exp(2\mathbf{s}^\top \mathbf{t})$ is a PDK
- ▶ $\exp(-\|\mathbf{s}\|^2) \exp(-\|\mathbf{t}\|^2)$ is a PDK as a product kernel
- ▶ the product of two PDK is a PDK

some examples of PD kernels...

type	name	$k(s, t)$
radial	gaussian	$\exp\left(-\frac{r^2}{b}\right), \quad r = \ s - t\ $
radial	laplacian	$\exp(-r/b)$
radial	rational	$1 - \frac{r^2}{r^2+b}$
radial	loc. gauss.	$\max\left(0, 1 - \frac{r}{3b}\right)^d \exp\left(-\frac{r^2}{b}\right)$
non stat.	χ^2	$\exp(-r/b), \quad r = \sum_k \frac{(s_k - t_k)^2}{s_k + t_k}$
projective	polynomial	$(s^\top t)^p$
projective	affine	$(s^\top t + b)^p$
projective	cosine	$s^\top t / \ s\ \ t\ $
projective	correlation	$\exp\left(\frac{s^\top t}{\ s\ \ t\ } - b\right)$

Most of the kernels depends on a quantity b called the bandwidth

kernels for objects and structures

kernels on histograms and probability distributions

$$k(p(x), q(x)) = \int k_i(p(x), q(x)) \mathbb{P}(x) dx$$

kernel on strings

- ▶ spectral string kernel
- ▶ using sub sequences
- ▶ similarities by alignments

$$k(\mathbf{s}, \mathbf{t}) = \sum_u \phi_u(\mathbf{s}) \phi_u(\mathbf{t})$$

$$k(\mathbf{s}, \mathbf{t}) = \sum_{\pi} \exp(\beta(\mathbf{s}, \mathbf{t}, \pi))$$

kernels on graphs

- ▶ the pseudo inverse of the (regularized) graph Laplacian

$$L = D - A \quad A \text{ is the adjacency matrix } D \text{ the degree matrix}$$

- ▶ diffusion kernels
- ▶ subgraph kernel convolution (using random walks)

$$\frac{1}{Z(b)} \exp^{bL}$$

and kernels on heterogeneous data (image), HMM, automata...

different point of view about kernels

kernel and scalar product

$$k(\mathbf{s}, \mathbf{t}) = \langle \phi(\mathbf{s}), \phi(\mathbf{t}) \rangle_{\mathcal{H}}$$

kernel and distance

$$d(\mathbf{s}, \mathbf{t})^2 = k(\mathbf{s}, \mathbf{s}) + k(\mathbf{t}, \mathbf{t}) - 2k(\mathbf{s}, \mathbf{t})$$

kernel and covariance: a positive matrix is a covariance matrix

$$\mathbb{P}(\mathbf{f}) = \frac{1}{Z} \exp\left(-\frac{1}{2}(\mathbf{f} - \mathbf{f}_0)^\top K^{-1}(\mathbf{f} - \mathbf{f}_0)\right)$$

if $\mathbf{f}_0 = 0$ and $\mathbf{f} = K\boldsymbol{\alpha}$, $\mathbb{P}(\boldsymbol{\alpha}) = \frac{1}{Z} \exp -\frac{1}{2}\boldsymbol{\alpha}^\top K\boldsymbol{\alpha}$

Kernel and regularity (green's function)

$$k(\mathbf{s}, \mathbf{t}) = P^* P \delta_{\mathbf{s}-\mathbf{t}} \quad \text{for some operator } P \quad (\text{e.g. some differential})$$

Roadmap

1 Kernels and the learning problem

- Two learning problems
- Kernelizing the linear regression
- Kernel machines: a definition

2 Tools: the functional framework

- In the beginning was the kernel
- Kernel and hypothesis set

3 Kernel machines

- Non sparse kernel machines
- sparse kernel machines: SVM
- practical SVM

4 Conclusion

From kernel to functions

$$\mathcal{H}_0 = \left\{ f \mid m_f < \infty; f_j \in \mathbb{R}; t_j \in \mathcal{X}, f(\mathbf{x}) = \sum_{j=1}^{m_f} f_j k(\mathbf{x}, t_j) \right\}$$

let define the bilinear form $(g(\mathbf{x}) = \sum_{i=1}^{m_g} g_i k(\mathbf{x}, s_i)) :$

$$\forall f, g \in \mathcal{H}_0, \langle f, g \rangle_{\mathcal{H}_0} = \sum_{j=1}^{m_f} \sum_{i=1}^{m_g} f_j g_i k(t_j, s_i)$$

Evaluation functional: $\forall \mathbf{x} \in \mathcal{X}$

$$f(\mathbf{x}) = \langle f(\cdot), k(\mathbf{x}, \cdot) \rangle_{\mathcal{H}_0}$$

from k to \mathcal{H}

with any positive kernel, a hypothesis set can be constructed \mathcal{H} with its metric

RKHS

Definition (reproducing kernel Hilbert space (RKHS))

a Hilbert space \mathcal{H} embedded with the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ is said to be with reproducing kernel if it exists a positive kernel k such that

$$\forall s \in \mathcal{X}, k(\cdot, s) \in \mathcal{H} \text{ et } \forall f \in \mathcal{H}, f(s) = \langle f(\cdot), k(s, \cdot) \rangle_{\mathcal{H}}$$

positive kernel \Leftrightarrow RKHS

- ▶ any function is pointwise defined
- ▶ defines the inner product
- ▶ it defines the **regularity** (smoothness) of the hypothesis set

functional differentiation in RKHS

Let J be a functional

$$J: \mathcal{H} \rightarrow \mathbb{R} \quad \text{examples:} \quad J_1(f) = \|f\|^2, J_2(f) = f(\mathbf{x}),$$
$$f \mapsto J(f)$$

J directional derivative in direction g at point f

$$dJ(f, g) = \lim_{\varepsilon \rightarrow 0} \frac{J(f + \varepsilon g) - J(f)}{\varepsilon}$$

Gradient $\nabla_J(f)$

$$\nabla_J: \mathcal{H} \rightarrow \mathcal{H} \quad \text{si} \quad dJ(f, g) = \langle \nabla_J(f), g \rangle_{\mathcal{H}}$$
$$f \mapsto \nabla_J(f)$$

exercice: find out $\nabla_{J_1}(f)$ et $\nabla_{J_2}(f)$

other kernels (what really matters)

- ▶ finite kernels

$$k(\mathbf{s}, \mathbf{t}) = (\phi_1(\mathbf{s}), \dots, \phi_p(\mathbf{s}))^\top (\phi_1(\mathbf{t}), \dots, \phi_p(\mathbf{t}))$$

- ▶ Mercer kernels

positive on a compact set

$$\Leftrightarrow k(\mathbf{s}, \mathbf{t}) = \sum_{j=1}^p \lambda_j \phi_j(\mathbf{s}) \phi_j(\mathbf{t})$$

- ▶ **positive kernels**

- ▶ positive semi-definite

- ▶ conditionally positive (for some functions p_j)

$$\forall \{\mathbf{x}_i\}_{i=1,n}, \forall \alpha_i, \sum_i^n \alpha_i p_j(\mathbf{x}_i) = 0; \quad j = 1, p, \quad \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 0$$

- ▶ symmetric non positive

$$k(\mathbf{s}, \mathbf{t}) = \tanh(\mathbf{s}^\top \mathbf{t} + \alpha_0)$$

- ▶ non symmetric – non positive

the key property: $\nabla_{J_t}(f) = k(t, \cdot)$ holds

Let's summarize

- ▶ positive kernels \Leftrightarrow RKHS $= \mathcal{H} \Leftrightarrow$ regularity $\|f\|_{\mathcal{H}}^2$
- ▶ the key property: $\nabla_{J_t}(f) = k(t, \cdot)$ holds not only for positive kernels
 $f(\mathbf{x}_i)$ exists (pointwise defined functions)
- ▶ universal consistency in RKHS
- ▶ the Gram matrix summarize the pairwise comparizons

Plan

- 1 Kernels and the learning problem**
 - Two learning problems
 - Kernelizing the linear regression
 - Kernel machines: a definition

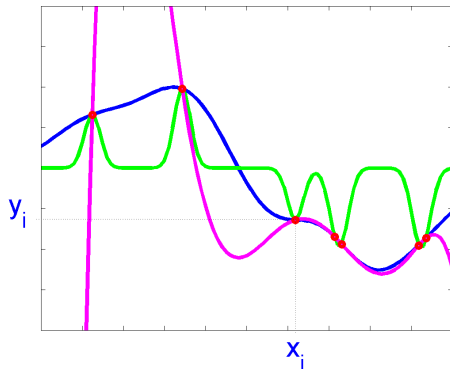
- 2 Tools: the functional framework**
 - In the beginning was the kernel
 - Kernel and hypothesis set

- 3 Kernel machines**
 - Non sparse kernel machines
 - sparse kernel machines: SVM
 - practical SVM

- 4 Conclusion**

Interpolation splines

find out $f \in \mathcal{H}$ such that $f(\mathbf{x}_i) = y_i$, $i = 1, \dots, n$



It is an ill posed problem

Interpolation splines: minimum norm interpolation

$$\left\{ \begin{array}{l} \min_{f \in \mathcal{H}} \quad \frac{1}{2} \|f\|_{\mathcal{H}}^2 \\ \text{such that} \quad f(\mathbf{x}_i) = y_i, \quad i = 1, \dots, n \end{array} \right.$$

The lagrangian (α_i Lagrange multipliers)

$$L(f, \alpha) = \frac{1}{2} \|f\|^2 - \sum_{i=1}^n \alpha_i (f(\mathbf{x}_i) - y_i)$$

optimality for f : $\nabla_f L(f, \alpha) = 0 \iff f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$

dual formulation (remove f from the lagrangian):

$$Q(\alpha) = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^n \alpha_i y_i \quad \text{solution: } \max_{\alpha \in \mathbb{R}^n} Q(\alpha)$$

$$K\alpha = y$$

Interpolation splines: minimum norm interpolation

$$\begin{cases} \min_{f \in \mathcal{H}} & \frac{1}{2} \|f\|_{\mathcal{H}}^2 \\ \text{such that} & f(\mathbf{x}_i) = y_i, \quad i = 1, \dots, n \end{cases}$$

The lagrangian (α_i Lagrange multipliers)

$$L(f, \alpha) = \frac{1}{2} \|f\|^2 - \sum_{i=1}^n \alpha_i (f(\mathbf{x}_i) - y_i)$$

optimality for f : $\nabla_f L(f, \alpha) = 0 \Leftrightarrow f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$

dual formulation (remove f from the lagrangian):

$$Q(\alpha) = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^n \alpha_i y_i \quad \text{solution: } \max_{\alpha \in \mathbb{R}^n} Q(\alpha)$$

$K\alpha = y$

Representer theorem

Theorem (epresenter theorem)

Let \mathcal{H} be a RKHS with kernel $k(s, t)$. Let ℓ be a function from \mathcal{X} to \mathbb{R} (loss function) and Φ a non decreasing function from \mathbb{R} to \mathbb{R} . If there exists a function f^* minimizing:

$$f^* = \operatorname{argmin}_{f \in \mathcal{H}} \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i)) + \Phi(\|f\|_{\mathcal{H}}^2)$$

then there exists a vector $\alpha \in \mathbb{R}^n$ such that:

$$f^*(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i)$$

it can be generalized to the semi parametric case: $+\sum_{j=1}^m \beta_j \phi_j(\mathbf{x})$

Smoothing splines

introducing the error (the slack) $\xi = f(x_i) - y_i$

$$(S) \quad \left\{ \begin{array}{l} \min_{f \in \mathcal{H}} \quad \frac{1}{2} \|f\|_{\mathcal{H}}^2 + \frac{1}{2\lambda} \sum_{i=1}^n \xi_i^2 \\ \text{such that} \quad f(x_i) = y_i + \xi_i, \quad i = 1, n \end{array} \right.$$

three equivalent definitions

$$(S') \quad \min_{f \in \mathcal{H}} \quad \frac{1}{2} \sum_{i=1}^n (f(x_i) - y_i)^2 + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2$$

$$\left\{ \begin{array}{l} \min_{f \in \mathcal{H}} \quad \frac{1}{2} \|f\|_{\mathcal{H}}^2 \\ \text{such that} \quad \sum_{i=1}^n (f(x_i) - y_i)^2 \leq C' \end{array} \right. \quad \left\{ \begin{array}{l} \min_{f \in \mathcal{H}} \quad \sum_{i=1}^n (f(x_i) - y_i)^2 \\ \text{such that} \quad \|f\|_{\mathcal{H}}^2 \leq C'' \end{array} \right.$$

using the representer theorem

$$(S'') \quad \min_{\alpha \in \mathbb{R}^n} \quad \frac{1}{2} \|K\alpha - \mathbf{y}\|^2 + \frac{\lambda}{2} \alpha^\top K \alpha$$

solution:

$$(S) \Leftrightarrow (S') \Leftrightarrow (S'') \Leftrightarrow (K + \lambda I)\alpha = \mathbf{y}$$

Kernel logistic regression

inspiration: the Bayes rule

$$D(\mathbf{x}) = \text{sign}(f(\mathbf{x}) + \alpha_0) \implies \log\left(\frac{\mathbb{P}(Y=1|\mathbf{x})}{\mathbb{P}(Y=-1|\mathbf{x})}\right) = f(\mathbf{x}) + \alpha_0$$

probabilities:

$$\mathbb{P}(Y = 1|\mathbf{x}) = \frac{\exp^{f(\mathbf{x})+\alpha_0}}{1 + \exp^{f(\mathbf{x})+\alpha_0}} \quad \mathbb{P}(Y = -1|\mathbf{x}) = \frac{1}{1 + \exp^{f(\mathbf{x})+\alpha_0}}$$

Rademacher distribution

$$\mathcal{L}(x_i, y_i, f, \alpha_0) = \mathbb{P}(Y = 1|\mathbf{x}_i)^{\frac{y_i+1}{2}} (1 - \mathbb{P}(Y = 1|\mathbf{x}_i))^{\frac{1-y_i}{2}}$$

penalized likelihood

$$\begin{aligned} J(f, \alpha_0) &= -\sum_{i=1}^n \log(\mathcal{L}(x_i, y_i, f, \alpha_0)) + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2 \\ &= \sum_{i=1}^n \log\left(1 + \exp^{-y_i(f(x_i) + \alpha_0)}\right) + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2 \end{aligned}$$

Let's summarize

- ▶ pros
 - ▶ Universality
 - ▶ from \mathcal{H} to \mathbb{R}^n using the representer theorem
 - ▶ no (explicit) curse of dimensionality

- ▶ splines $\mathcal{O}(n^3)$ (can be reduced to $\mathcal{O}(n^2)$)

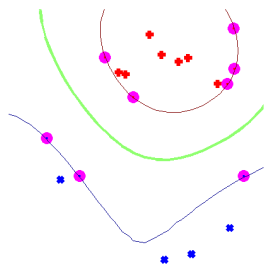
- ▶ logistic regression $\mathcal{O}(kn^3)$ (can be reduced to $\mathcal{O}(kn^2)$)

- ▶ no scalability!

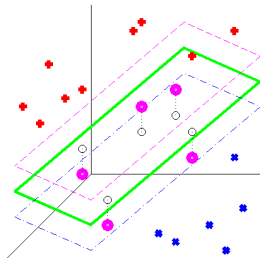
sparseness comes to the rescue!

using relevant features...

a data point becomes a function $\mathbf{x} \rightarrow k(\mathbf{x}, \cdot)$



input space representation: x



feature space: $k(x, \cdot)$

SVM dual formulation

Dual formulation

$$\left\{ \begin{array}{l} \max_{\alpha \in \mathbb{R}^n} \quad -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^n \alpha_i \\ \text{with} \quad \sum_{i=1}^n \alpha_i y_i = 0 \quad \text{and} \quad 0 \leq \alpha_i, \quad i = 1, n \end{array} \right.$$

The dual formulation gives a quadratic program (QP)

$$\left\{ \begin{array}{l} \min_{\alpha \in \mathbb{R}^n} \quad \frac{1}{2} \alpha^\top G \alpha - \mathbb{1}^\top \alpha \\ \text{with} \quad \alpha^\top \mathbf{y} = 0 \quad \text{and} \quad 0 \leq \alpha \end{array} \right.$$

with $G_{ij} = y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$

with the linear kernel $f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i (\mathbf{x}^\top \mathbf{x}_i) = \sum_{j=1}^d \beta_j \mathbf{x}_j$
when d is small wrt. n primal may be interesting.

the general case: C-SVM

Primal formulation

$$(\mathcal{P}) \begin{cases} \min_{f \in \mathcal{H}, \alpha_0, \xi \in \mathbb{R}^n} & \frac{1}{2} \|f\|^2 + \frac{C}{p} \sum_{i=1}^n \xi_i^p \\ \text{such that} & y_i (f(\mathbf{x}_i) + \alpha_0) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, n \end{cases}$$

C is the *regularization path* parameter (to be tuned)

$p = 1, L_1$ SVM

$$\begin{cases} \max_{\alpha \in \mathbb{R}^n} & -\frac{1}{2} \alpha^\top H \alpha + \alpha^\top \mathbb{I} \\ \text{such that} & \alpha^\top \mathbf{y} = 0 \text{ and } 0 \leq \alpha_i \leq C \quad i = 1, n \end{cases}$$

$p = 2, L_2$ SVM

$$\begin{cases} \max_{\alpha \in \mathbb{R}^n} & -\frac{1}{2} \alpha^\top (H + \frac{1}{C} I) \alpha + \alpha^\top \mathbb{I} \\ \text{such that} & \alpha^\top \mathbf{y} = 0 \text{ and } 0 \leq \alpha_i \quad i = 1, n \end{cases}$$

the regularization path: is the set of solutions $\alpha(C)$ when C varies

The importance of being support

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i k(\mathbf{x}_i, \mathbf{x})$$

data point	α	constraint value
\mathbf{x}_i useless	$\alpha_i = 0$	$y_i (f(\mathbf{x}_i) + \alpha_0) > 1$
\mathbf{x}_i support	$\alpha_i > 0$	$y_i (f(\mathbf{x}_i) + \alpha_0) = 1$

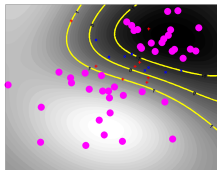
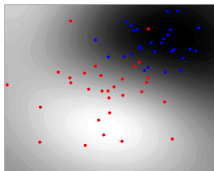
Table: When a data point is « support » it lies exacty on the margin.

here lies the efficiency of the algorithm (and its complexity)!

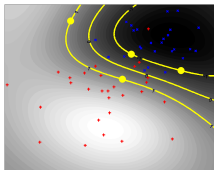
sparsness: $\alpha_j = 0$

Data groups: illustration

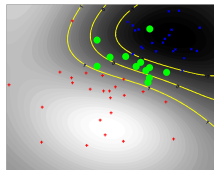
$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i) + \alpha_0$$
$$D(x) = \text{sign}(f(\mathbf{x}))$$



useless data
well classified
 $\alpha = 0$



important data
support
 $0 < \alpha < C$



suspicious data
 $\alpha = C$

Two more ways to derivate SVM

Using the hinge loss

$$\min_{f \in \mathcal{H}, \alpha_0 \in \mathbb{R}} \frac{1}{p} \sum_{i=1}^n \max(0, 1 - y_i(f(\mathbf{x}_i) + \alpha_0))^p + \frac{1}{2C} \|f\|_{\mathcal{H}}^2$$

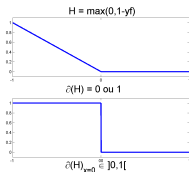
Minimizing the distance between the convex hulls

$$\left\{ \begin{array}{l} \min_{\alpha} \quad \|u - v\|_{\mathcal{H}}^2 \\ \text{with} \quad u(\mathbf{x}) = \sum_{\{i|y_i=1\}} \alpha_i k(\mathbf{x}_i, \mathbf{x}), \quad v(\mathbf{x}) = \sum_{\{i|y_i=-1\}} \alpha_i k(\mathbf{x}_i, \mathbf{x}) \\ \text{and} \quad \sum_{\{i|y_i=1\}} \alpha_i = 1, \quad \sum_{\{i|y_i=-1\}} \alpha_i = 1, \quad 0 \leq \alpha_i \quad i = 1, n \end{array} \right.$$

$$f(\mathbf{x}) = \frac{2}{\|u - v\|_{\mathcal{H}}^2} (u(\mathbf{x}) - v(\mathbf{x})) \quad \text{and} \quad \alpha_0 = \frac{\|u\|_{\mathcal{H}}^2 - \|v\|_{\mathcal{H}}^2}{\|u - v\|_{\mathcal{H}}^2}$$

Regularization path for SVM

$$\min_{f \in \mathcal{H}} \sum_{i=1}^n \max(1 - y_i f(\mathbf{x}_i), 0) + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2$$



I_α is the set of support vectors s.t. $y_i f(\mathbf{x}_i) = 1$;

$$\nabla_f J(f) = \sum_{i \in I_\alpha} \alpha_i y_i K(\mathbf{x}_i, \cdot) + \sum_{i \in \mathcal{H}} y_i K(\mathbf{x}_i, \cdot) + \lambda f(\cdot) \quad \text{with} \quad \alpha_i = \partial H(\mathbf{x}_i)$$

in particular at point $\mathbf{x}_j \in I_\alpha$ ($f_0(\mathbf{x}_j) = f_n(\mathbf{x}_j) = y_j$)

$$\begin{aligned} \sum_{i \in I_\alpha} \alpha_{i0} y_i K(\mathbf{x}_i, \mathbf{x}_j) &= - \sum_{i \in \mathcal{H}} y_i K(\mathbf{x}_i, \mathbf{x}_j) - \lambda_0 y_j \\ \sum_{i \in I_\alpha} \alpha_{in} y_i K(\mathbf{x}_i, \mathbf{x}_j) &= - \sum_{i \in \mathcal{H}} y_i K(\mathbf{x}_i, \mathbf{x}_j) - \lambda_n y_j \end{aligned}$$

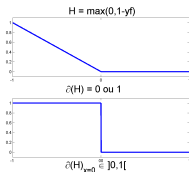
$$G(\alpha_n - \alpha_0) = (\lambda_0 - \lambda_n) y \quad \text{avec} \quad G_{ij} = y_i K(\mathbf{x}_i, \mathbf{x}_j)$$

$$\alpha_n = \alpha_0 + (\lambda_0 - \lambda_n) w$$

$$w = (G)^{-1} y$$

Regularization path for SVM

$$\min_{f \in \mathcal{H}} \sum_{i=1}^n \max(1 - y_i f(\mathbf{x}_i), 0) + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2$$



I_α is the set of support vectors s.t. $y_i f(\mathbf{x}_i) = 1$;

$$\nabla_f J(f) = \sum_{i \in I_\alpha} \alpha_i y_i K(\mathbf{x}_i, \cdot) + \sum_{i \in \mathcal{H}} y_i K(\mathbf{x}_i, \cdot) + \lambda f(\cdot) \quad \text{with} \quad \alpha_i = \partial H(\mathbf{x}_i)$$

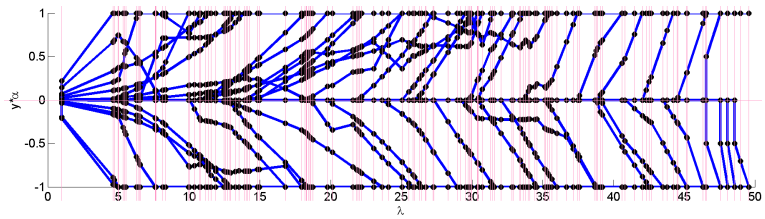
in particular at point $\mathbf{x}_j \in I_\alpha$ ($f_o(\mathbf{x}_j) = f_n(\mathbf{x}_j) = y_j$)

$$\frac{\sum_{i \in I_\alpha} \alpha_{i0} y_i K(\mathbf{x}_i, \mathbf{x}_j) = - \sum_{i \in \mathcal{H}} y_i K(\mathbf{x}_i, \mathbf{x}_j) - \lambda_o y_j}{\sum_{i \in I_\alpha} \alpha_{in} y_i K(\mathbf{x}_i, \mathbf{x}_j) = - \sum_{i \in \mathcal{H}} y_i K(\mathbf{x}_i, \mathbf{x}_j) - \lambda_n y_j} = \frac{(\lambda_o - \lambda_n) \mathbf{y}}{G} \quad \text{avec} \quad G_{ij} = y_i K(\mathbf{x}_i, \mathbf{x}_j)$$

$$\alpha_n = \alpha_o + (\lambda_o - \lambda_n) \mathbf{w}$$

$$\mathbf{w} = (G)^{-1} \mathbf{y}$$

Example of regularization path



α estimation and data selection

How to choose ℓ and P to get linear regularization path?

the *path* is piecewise linear \Leftrightarrow one is piecewise quadratic and the other is piecewise linear

the convex case [Rosset & Zhu, 07]

$$\min_{\beta \in \mathbb{R}^d} \ell(\beta) + \lambda P(\beta)$$

1. piecewise linearity: $\lim_{\varepsilon \rightarrow 0} \frac{\beta(\lambda + \varepsilon) - \beta(\lambda)}{\varepsilon} = \text{constant}$

2. optimality

$$\begin{aligned} \nabla \ell(\beta(\lambda)) + \lambda \nabla P(\beta(\lambda)) &= 0 \\ \nabla \ell(\beta(\lambda + \varepsilon)) + (\lambda + \varepsilon) \nabla P(\beta(\lambda + \varepsilon)) &= 0 \end{aligned}$$

3. Taylor expansion

$$\lim_{\varepsilon \rightarrow 0} \frac{\beta(\lambda + \varepsilon) - \beta(\lambda)}{\varepsilon} = [\nabla^2 \ell(\beta(\lambda)) + \lambda \nabla^2 P(\beta(\lambda))]^{-1} \nabla P(\beta(\lambda))$$

$\nabla^2 \ell(\beta(\lambda)) = \text{constant}$ and $\nabla^2 P(\beta(\lambda)) = 0$

Problems with Piecewise linear regularization path

L	P	<i>regression</i>	<i>classification</i>	<i>clustering</i>
L_2	L_1	Lasso/LARS	L1 L2 SVM	PCA L1
L_1	L_2	SVR	SVM	OC SVM
L_1	L_1	L1 LAD Danzig Selector	L1 SVM	

Table: example of piecewise linear regularization path algorithms.

$$P : L_p = \sum_{j=1}^d |\beta_j|^p$$

$$L : L_p : |f(\mathbf{x}) - y|^p \quad \text{hinge } (yf(\mathbf{x}) - 1)_+^p$$

$$\varepsilon\text{-insensitive} \quad \begin{cases} 0 & \text{if } |f(\mathbf{x}) - y| < \varepsilon \\ |f(\mathbf{x}) - y| - \varepsilon & \text{else} \end{cases}$$

$$\text{Huber's loss:} \quad \begin{cases} |f(\mathbf{x}) - y|^2 & \text{if } |f(\mathbf{x}) - y| < t \\ 2t|f(\mathbf{x}) - y| - t^2 & \text{else} \end{cases}$$

standart formulation

- ▶ portfolio optimization (Markovitz, 1952)
 - ▶ return vs. risk

$$\begin{cases} \min_{\beta} & \frac{1}{2}\beta^T Q \beta \\ \text{with} & \mathbf{e}^T \beta = C \end{cases}$$



- ▶ *efficiency frontier*: piecewise linear (*Critical path Algo.*)

- ▶ Sensitivity analysis: standart formulation (Heller, 1954)

$$\begin{cases} \min_{\beta} & \frac{1}{2}\beta^T Q \beta + (\mathbf{c} + \lambda \Delta \mathbf{c})^T \beta \\ \text{with} & A\beta = \mathbf{b} + \mu \Delta \mathbf{b} \end{cases}$$

- ▶ Parametric programming (see T. Gal's book 1968)
 - ▶ in the general case of PLP: the reg. path is piecewise linear
 - ▶ ... and PQP is piecewise quadratic
 - ▶ multiparametric programming

ν-SVM and other formulations...

$\nu \in [0, 1]$

$$(\nu) \left\{ \begin{array}{l} \min_{f, \alpha_0, \xi, m} \quad \frac{1}{2} \|f\|_{\mathcal{H}}^2 + \frac{1}{np} \sum_{i=1}^n \xi_i^p - \nu m \\ \text{with} \quad y_i (f(\mathbf{x}_i) + \alpha_0) \geq m - \xi_i, \quad i = 1, n, \\ \text{and} \quad m \geq 0, \quad \xi_i \geq 0, \quad i = 1, n, \end{array} \right.$$

for $p = 1$ the dual formulation is:

$$\left\{ \begin{array}{l} \max_{\alpha \in \mathbf{R}^n} \quad -\frac{1}{2} \alpha^\top G \alpha \\ \text{with} \quad \alpha^\top \mathbf{y} = 0 \text{ et } 0 \leq \alpha_i \leq \frac{1}{n} \quad i = 1, n \\ \text{and} \quad \nu \leq \alpha^\top \mathbf{1} \end{array} \right.$$

$$C = \frac{1}{m}$$

SVM with non symmetric costs

problem in the primal

$$\left\{ \begin{array}{l} \min_{f \in \mathcal{H}, \alpha_0, \xi \in \mathbb{R}^n} \quad \frac{1}{2} \|f\|_{\mathcal{H}}^2 + C^+ \sum_{\{i|y_i=1\}} \xi_i^p + C^- \sum_{\{i|y_i=-1\}} \xi_i^p \\ \text{with} \quad y_i (f(\mathbf{x}_i) + \alpha_0) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, n \end{array} \right.$$

for $p = 1$ the dual formulation is the following:

$$\left\{ \begin{array}{l} \max_{\alpha \in \mathbb{R}^n} \quad -\frac{1}{2} \alpha^\top G \alpha + \alpha^\top \mathbb{I} \\ \text{with} \quad \alpha^\top \mathbf{y} = 0 \text{ and } 0 \leq \alpha_i \leq C^+ \text{ or } C^- \quad i = 1, n \end{array} \right.$$

Generalized SVM

$$\min_{f \in \mathcal{H}, \alpha_0 \in \mathbb{R}} \sum_{i=1}^n \max(0, 1 - y_i(f(\mathbf{x}_i) + \alpha_0)) + \frac{1}{C} \varphi(f) \quad \varphi \text{ convex}$$

in particular $\varphi(f) = \|f\|_p^p$ with $p = 1$ leads to L1 SVM.

$$\left\{ \begin{array}{l} \min_{\alpha \in \mathbb{R}^n, \alpha_0, \xi} \quad \mathbb{1}^\top \beta + C \mathbb{1}^\top \xi \\ \text{with} \quad y_i \left(\sum_{j=1}^n \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) + \alpha_0 \right) \geq 1 - \xi_i, \\ \text{and} \quad -\beta_i \leq \alpha_i \leq \beta_i, \quad \xi_i \geq 0, \quad i = 1, n \end{array} \right.$$

with $\beta = |\alpha|$. the dual is:

$$\left\{ \begin{array}{l} \max_{\gamma, \delta, \delta^* \in \mathbb{R}^{3n}} \quad \mathbb{1}^\top \gamma \\ \text{with} \quad \mathbf{y}^\top \gamma = 0, \quad \delta_i + \delta_i^* = 1 \\ \quad \sum_{j=1}^n \gamma_j k(\mathbf{x}_i, \mathbf{x}_j) = \delta_i - \delta_i^*, \quad i = 1, n \\ \text{and} \quad 0 \leq \delta_i, 0 \leq \delta_i^*, 0 \leq \gamma_i \leq C, \quad i = 1, n \end{array} \right.$$

SVM reduction (reduced set method)

- ▶ objective: compile the model
- ▶ $f(x) = \sum_{i=1}^{n_s} \alpha_i k(\mathbf{x}_i, \mathbf{x})$, $n_s \ll n$, n_s too big

▶ compiled model as the solution of:

$$g(\mathbf{x}) = \sum_{i=1}^{n_c} \beta_i k(\mathbf{z}_i, \mathbf{x}), n_c \ll n_s$$

- ▶ β, \mathbf{z}_i and c are tuned by minimizing:

$$\min_{\beta, \mathbf{z}_i} \|g - f\|_H^2$$

where

$$\min_{\beta, \mathbf{z}_i} \|g - f\|_H^2 = \alpha^\top K_x \alpha + \beta^\top K_z \beta - 2\alpha^\top K_{xz} \beta$$

some authors advice $0,03 \leq \frac{n_c}{n_s} \leq 0,1$

- ▶ solve it by using use (stochastic) gradient (its a RBF problem)

SVM and probabilities (1/2)

$\log \frac{\mathbb{P}(Y = 1|\mathbf{x})}{\mathbb{P}(Y = -1|\mathbf{x})}$ as (almost) the same sign as $f(\mathbf{x})$

$$\log \frac{\mathbb{P}(Y = 1|\mathbf{x})}{\mathbb{P}(Y = -1|\mathbf{x})} = a_1 f(\mathbf{x}) + a_2 \quad \mathbb{P}(Y = 1|\mathbf{x}) = \frac{1}{1 + \exp^{a_1 f(\mathbf{x}) + a_2}}$$

a_1 et a_2 estimated using maximum likelihood

some facts

- ▶ SVM is universally consistent (converges towards the Bayes risk)
- ▶ SVM asymptotically implements the bayes rule
- ▶ but theoretically: **no consistency towards conditional probabilities** (due to the nature of sparsity)
- ▶ to estimate conditional probabilities on an interval (typically $[\frac{1}{2} - \eta, \frac{1}{2} + \eta]$) to sparseness in this interval (all data points have to be support vectors)

SVM and probabilities (2/2)

An alternative approach

$$g(\mathbf{x}) - \varepsilon^-(\mathbf{x}) \leq \mathbb{P}(Y = 1|\mathbf{x}) \leq g(\mathbf{x}) + \varepsilon^+(\mathbf{x})$$

with $g(\mathbf{x}) = \frac{1}{1+4^{-f(\mathbf{x})-\alpha_0}}$

non parametric functions ε^- and ε^+ have to verify:

$$\begin{aligned} g(\mathbf{x}) + \varepsilon^+(\mathbf{x}) &= \exp^{-a_1(1-f(\mathbf{x})-\alpha_0)+a_2} \\ 1 - g(\mathbf{x}) - \varepsilon^-(\mathbf{x}) &= \exp^{-a_1(1+f(\mathbf{x})+\alpha_0)+a_2} \end{aligned}$$

with $a_1 = \log 2$ and $a_2 = 0$

logistic regression and the import vector machine

- ▶ Logistic regression is NON sparse
- ▶ kernalize it using the *dictionary* strategy
- ▶ Algorithm:
 - ▶ find the solution of the KLR using only a subset \mathcal{S} of the data
 - ▶ build \mathcal{S} iteratively using active constraint approach
- ▶ this trick brings sparsity
- ▶ it estimates probability
- ▶ it can naturally be generalized to the multiclass case

- ▶ efficient when uses:
 - ▶ a few import vectors
 - ▶ component-wise update procedure

- ▶ extention using L_1 KLR

Multiclass SVM

- ▶ one vs all: winner takes all
- ▶ one vs one:
 - ▶ max-wins voting
 - ▶ pairwise coupling: use probability
- ▶ global approach (size $c \times n$),
 - ▶ formal (differs variations)

$$\left\{ \begin{array}{l} \min_{f \in \mathcal{H}, \alpha_0, \xi \in \mathbb{R}^n} \frac{1}{2} \sum_{\ell=1}^c \|f_{\ell}\|_{\mathcal{H}}^2 + \frac{C}{p} \sum_{i=1}^n \sum_{\ell=1, \ell \neq y_i}^c \xi_{i\ell}^p \\ \text{with } y_i (f_{y_i}(\mathbf{x}_i) + b_{y_i} - f_{\ell}(\mathbf{x}_i) - b_{\ell}) \geq 1 - \xi_{i\ell} \\ \text{and } \xi_{i\ell} \geq 0 \text{ for } i = 1, \dots, n; \ell = 1, \dots, c; \ell \neq y_i \end{array} \right.$$

non consistent estimator but practically useful

- ▶ structured outputs

approach	problem size	number of sub problems
<i>all together</i>	$n \times c$	1
<i>1 vs. all</i>	n	c
<i>1 vs. 1</i>	$\frac{2n}{c}$	$\frac{c(c-1)}{2}$

Multiclass SVM

- ▶ one vs all: winner takes all
- ▶ one vs one:
 - ▶ max-wins voting
 - ▶ pairwise coupling: use probability – best results
- ▶ global approach (size $c \times n$),
 - ▶ formal (differs variations)

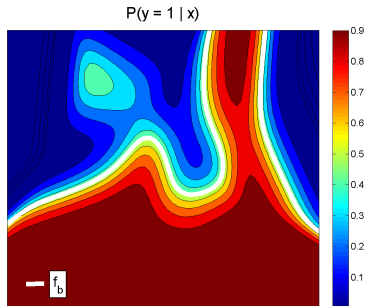
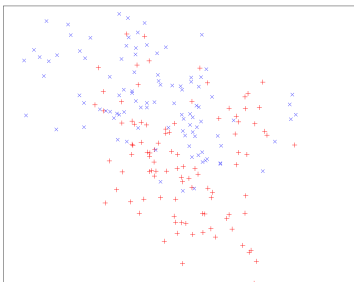
$$\left\{ \begin{array}{l} \min_{f \in \mathcal{H}, \alpha_0, \xi \in \mathbb{R}^n} \quad \frac{1}{2} \sum_{\ell=1}^c \|f_{\ell}\|_{\mathcal{H}}^2 + \frac{C}{p} \sum_{i=1}^n \sum_{\ell=1, \ell \neq y_i}^c \xi_{i\ell}^p \\ \text{with } y_i (f_{y_i}(\mathbf{x}_i) + b_{y_i} - f_{\ell}(\mathbf{x}_i) - b_{\ell}) \geq 1 - \xi_{i\ell} \\ \text{and } \xi_{i\ell} \geq 0 \text{ for } i = 1, \dots, n; \ell = 1, \dots, c; \ell \neq y_i \end{array} \right.$$

non consistent estimator but practicaly useful

- ▶ structured outputs

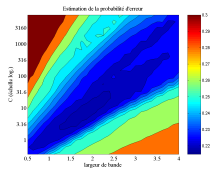
approach	problem size	number of sub problems
<i>all together</i>	$n \times c$	1
<i>1 vs. all</i>	n	c
<i>1 vs. 1</i>	$\frac{2n}{c}$	$\frac{c(c-1)}{2}$

Mixture data

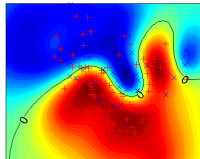


- ▶ $x : 200 \times 2$
- ▶ $y : 100$ for each class
- ▶ mixture model with 10 gaussians
- ▶ the bayes error is known

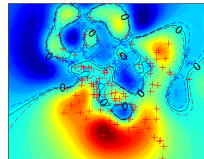
C and σ influence



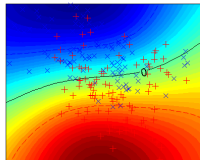
C = 1; $\sigma = 1$



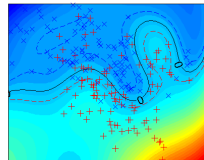
C = 10000; $\sigma = 1$



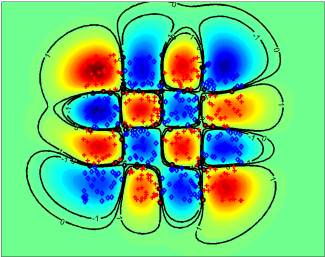
C = 1; $\sigma = 5$



C = 10000; $\sigma = 5$

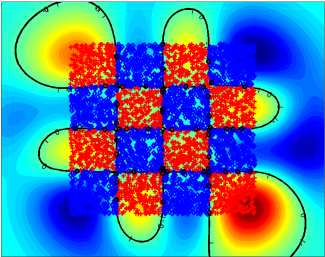


a separable case

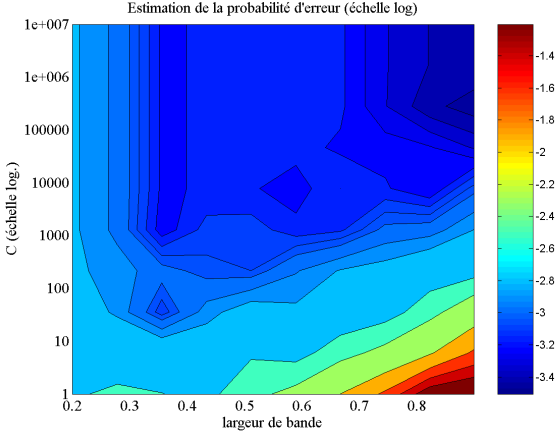


$n = 500$ data points

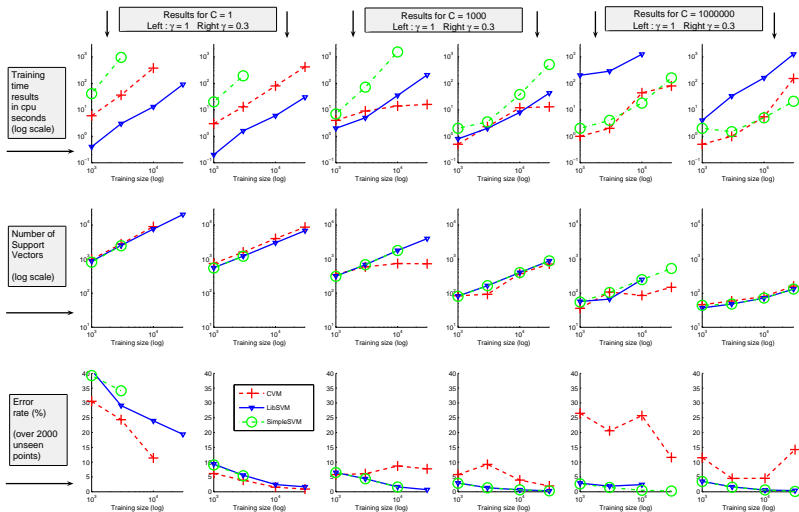
$n = 5000$ data points



tuning C and σ : grid search



empirical complexity



Conclusion

- ▶ nonlinearity through kernel: using examples influence
- ▶ universality: kernel to functions (R.K.H.S.)
- ▶ representer theorem: from functions to vectors
- ▶ $L1$ provides sparsity

a question of vocabulary

- ▶ margin: **regularization**
- ▶ Mercer kernel: **positive kernel**
- ▶ SVM: **a method among others**

- ▶ kernels (RKHS) universality
- ▶ regularization univ. consistency
- ▶ convexity efficiency
- ▶ sparsity efficiency

no (explicit) model

but a kernel, a cost and a regularity

challenges: towards tough learning

- ▶ the size effect
 - ▶ ready to use: automatization
 - ▶ adaptative: on line context aware
 - ▶ beyond kenrels
- ▶ Automatic and adaptive model selection
 - ▶ variable selection
 - ▶ kernel tuning (k et σ)
 - ▶ hyperparametres: C , duality gap, λ
- ▶ \mathbb{P} change
- ▶ Theory
 - ▶ non positive kernels
 - ▶ a more general representer theorem

biblio: kernel-machines.org

- ▶ John Shawe-Taylor and Nello Cristianini *Kernel Methods for Pattern Analysis*, Cambridge University Press, 2004
- ▶ Bernhard Schölkopf and Alex Smola. *Learning with Kernels*. MIT Press, Cambridge, MA, 2002.
- ▶ Trevor Hastie, Robert Tibshirani and Jerome Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, Springer, 2001
- ▶ Léon Bottou, Olivier Chapelle, Dennis DeCoste and Jason Weston *Large-Scale Kernel Machines (Neural Information Processing)*, MIT press 2007
- ▶ Olivier Chapelle, Bernhard Scholkopf and Alexander Zien, *Semi-supervised Learning*, MIT press 2006
- ▶ Vladimir Vapnik. *Estimation of Dependences Based on Empirical Data*. Springer Verlag, 2006, 2nd edition.
- ▶ Vladimir Vapnik. *The Nature of Statistical Learning Theory*. Springer, 1995.
- ▶ Grace Wahba. *Spline Models for Observational Data*. SIAM CBMS-NSF Regional Conference Series in Applied Mathematics vol. 59, Philadelphia, 1990
- ▶ Alain Berlinet and Christine Thomas-Agnan, *Reproducing Kernel Hilbert Spaces in Probability and Statistics*, Kluwer Academic Publishers, 2003
- ▶ Marc Atteia et Jean Gaches , *Approximation Hilbertienne - Splines, Ondelettes, Fractales*, PUG, 1999