

Hilbert Space Methods in Learning

guest lecturer: Risi Kondor

6772 Advanced Machine Learning and Perception (Jebara),
Columbia University, October 15, 2003.

1. A general formulation of the learning problem

- Empirical and true errors — overfitting
- Error bounds and what they tell us about the design of algorithms

2. Hilbert space methods

- Reproducing Kernel Hilbert Spaces
- Kernels
- Algorithms: SVM, Gaussian Processes, Kernel PCA

Tutorial online at

<http://www.cs.columbia.edu/risi/notes/tutorial6672.pdf>

The Learning Problem

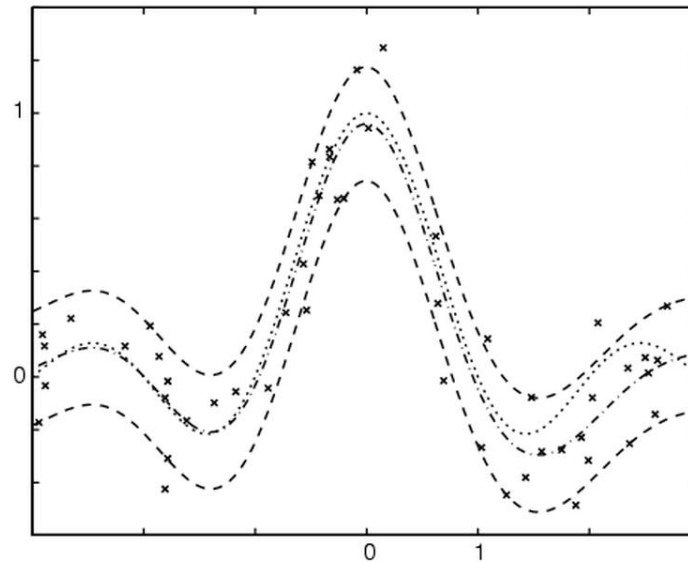
Regression

Learn a function $f : x \mapsto y$

Linear functions, order p polynomials, splines, etc.

Examples:

Boston housing problem, robot grasps, motorcycle data, etc.

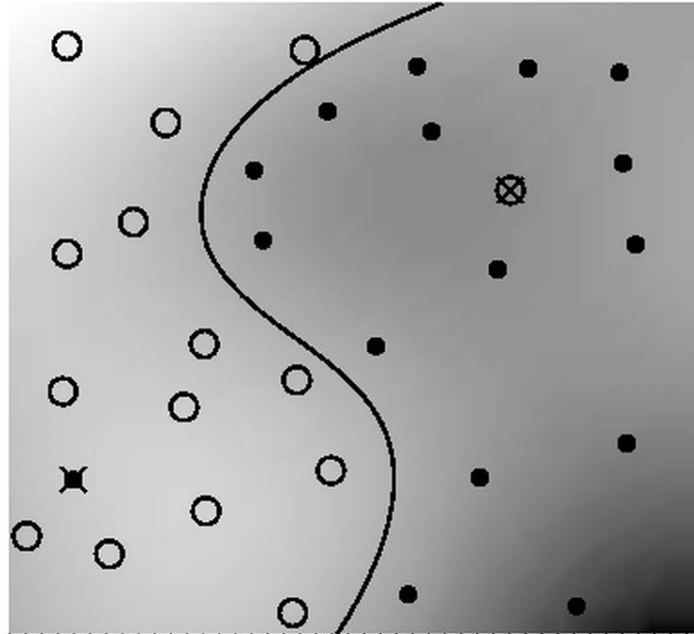


Classification

Separate $+1$ labeled points from -1 labeled points

Examples:

Face recognition, DNA splice site identification, document classification, call type classification



Supervised learning

Input space: \mathcal{X}

e.g. $\mathcal{X} = \mathbb{R}^n$

Output space: \mathcal{Y}

$\mathcal{Y} = \{-1, +1\}$ for classification

$\mathcal{Y} = \mathbb{R}$ for regression

Training set: $S = (x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$

$x_i \in \mathcal{X}, y_t \in \mathcal{Y}$

“Truth”:

- Deterministic: $y = f_0(x)$
- Probabilistic: $y \sim p(y|x)$ (more general)

Goal: construct **hypothesis** $f : \mathcal{X} \mapsto \mathcal{Y}$ to predict y given x .

The Empirical Risk

Empirical risk (training error):

$$R_{\text{emp}}[f] = \frac{1}{m} \sum_{i=1}^m L(f(x_i), y_i)$$

where $L : \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}$ is the **loss function**.

Zero-one loss for classification: $L(\hat{y}, y) = \begin{cases} 1 & \text{if } \hat{y} \neq y \\ 0 & \text{otherwise.} \end{cases}$

Squared error loss for regression: $L(\hat{y}, y) = (y - \hat{y})^2$.

A Bad Learning Algorithm (memorization algorithm)

Set

$$f(x) = \begin{cases} 1 & \text{when } x = x_i \text{ and } y_i = 1 \\ -1 & \text{otherwise.} \end{cases}$$

For zero-one loss perfect performance on training data!

$$R_{\text{emp}}[f] = \frac{1}{m} \sum_{i=1}^m L(f(x_i), y_i) = 0$$

Will it **generalize** well to testing examples? Why not?

The True Risk

- Assume some distribution on inputs: $p(x)$
- Distribution on (x, y) examples: $p(x, y) = p(y | x) p(x)$ or $p(x, y) = \delta(y - f_0(x)) p(x)$

True risk:

$$R[f] = \mathbb{E} [L(f(x), y)] = \int_{\mathcal{X} \times \mathcal{Y}} p(x, y) L(f(x), y) dx dy.$$

This is what we really want to minimize in **discriminative** learning.

True Risk vs. Empirical Risk

$$R[f] = \mathbb{E} [L(f(x), y)] \quad \longleftrightarrow \quad R_{\text{emp}}[f] = \frac{1}{m} \sum_{i=1}^m L(f(x_i), y_i).$$

Just minimizing R_{emp} is BAD (see previous algorithm). Optimizing the training error at the expense of the testing error is called **overfitting**.

But we do not know $p(x, y)$!!! Can we still do anything?

Bounding the True Risk

For many practical learning algorithms

$$R[f] = \mathbb{E} [L(f(x), y)] \quad \sim \quad R_{\text{emp}}[f] = \frac{1}{m} \sum_{i=1}^m L(f(x_i), y_i).$$

Uniform error bounds:

For any distribution D , with probability $1 - \delta$ (over the choice of training set)

$$| R[f] - R_{\text{emp}}[f] | \leq \epsilon$$

for all hypotheses $f \in \mathcal{F}$ simultaneously.

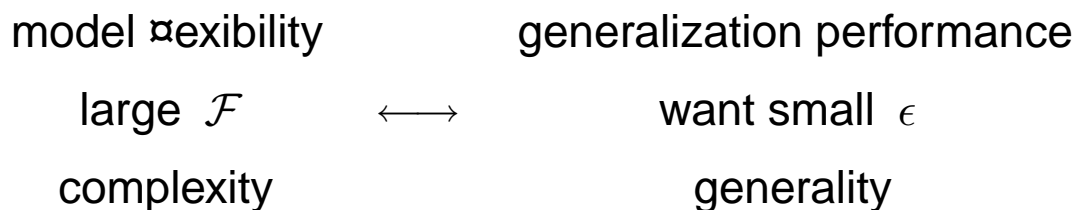
“PAC” bound: probably approximately correct

KEY CONCEPT: Capacity Control

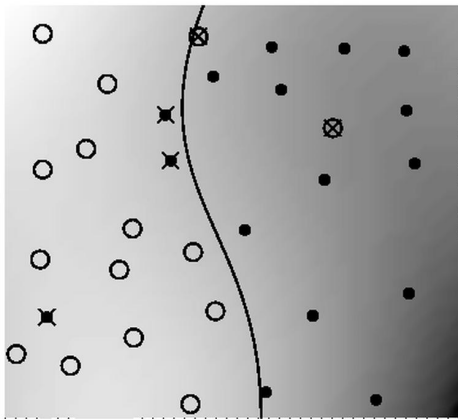
$$\mathbb{P} \left[|R[f] - R_{\text{emp}}[f]| \leq \epsilon \right] \geq 1 - \delta$$

Generally, ϵ is a complicated function of δ depending crucially on \mathcal{F} (**hypothesis class**) that f is chosen from.

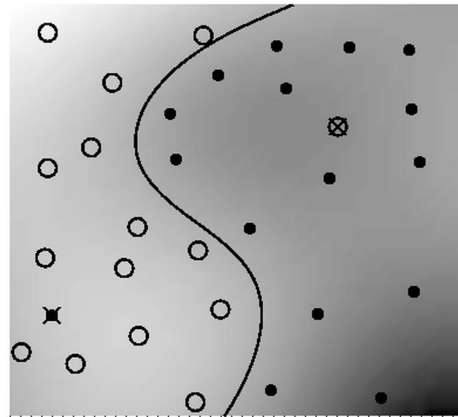
Compromise:



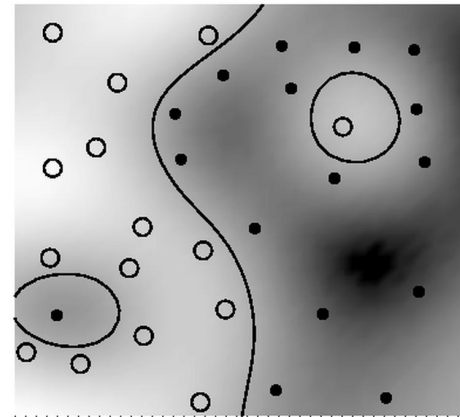
Capacity control



Too inflexible?



Just right.



Overfitting?

How do we quantify complexity of f ?

Uniform Error Bounds

$$\mathbb{P} \left[R[f] - R_{\text{emp}}[f] \leq \epsilon \mid \forall f \in \mathcal{F} \right] \geq 1 - \delta.$$

$$\mathbb{P} \left[\sup_{f \in \mathcal{F}} [R[f] - R_{\text{emp}}[f]] \leq \epsilon \right] \geq 1 - \delta,$$

Not equivalent to:

$$\mathbb{P} \left[R[f] - R_{\text{emp}}[f] \leq \epsilon \right] \geq 1 - \delta \quad \forall f \in \mathcal{F}.$$

Vapnik-Chervonenkis type bounds

With probability $1 - \delta$

$$\sup_{f \in \mathcal{F}} [R[f] - R_{\text{emp}}[f]] \leq \sqrt{\frac{h (\log(2m/h) + 1) - \log(\delta/4)}{m}}$$

where h is the VC dimension of \mathcal{F} .

Linear discriminators in \mathbb{R}^n : $h = n + 1$

... with margin γ in ball of radius D : $h = \min(n, \lceil D^2/\gamma^2 \rceil) + 1$

Large margin is good!

Covering number bounds

With probability $1 - \delta$

$$\sup_{f \in \mathcal{F}} [R[f] - R_{\text{emp}}[f]] \leq 16M \sqrt{\frac{\log(12m \mathbb{E} \mathcal{N}_1(S, \epsilon/8)) - \log \delta}{m}}$$

where M is an upper bound on $|L(f(x), y)|$.

The covering number $\mathcal{N}_1(S, \epsilon)$ is the number of vectors v_1, v_2, \dots, v_n needed to ensure that for any $f \in \mathcal{F}$, there is a v_k such that

$$\sum_{i=1}^m L(f(x_i), v_k) \leq \epsilon.$$

Stability-based bounds

If f^* is the hypothesis returned by a β -stable algorithm, then with probability $1 - \delta$

$$R[f^*] - R_{\text{emp}}[f^*] \leq \beta + \sqrt{\frac{-2(m\beta + M)^2 \log(\delta/2)}{m}}$$

where M is an upper bound on $|L(f(x), y)|$.

An algorithm is β -stable if for all training sets, and any example (x, y) , $L(f^*(x), y)$ changes by at most β when we replace any one of the training examples by any other example.

Rademacher bounds

For $\text{Ra}_r[f] = r$ with probability $1 - \delta$

$$R[f] \leq \inf_{\alpha > 0} \left[(1 + \alpha) R_{\text{emp}}[f] + \left(1 + \frac{1}{4\alpha}\right) \left(31r \log^2 \frac{b}{r} + 50 \frac{b\epsilon}{n}\right) \right].$$

Rademacher average:

$$\text{Ra}_r[f] = \mathbb{E}_{S, \sigma} \left[\sup_{f \in \mathcal{F} : \mathbb{E}L(f(x), y) \leq r} \sum_{i=1}^m \sigma_i L(f(x_i), y_i) \right]$$

where $\mathbb{P}[\sigma = 1] = \mathbb{P}[\sigma = -1] = 1/2$.

Structural Risk Minimization

If we have bound of form

$$\mathbb{P} \left[\sup_{f \in \mathcal{F}} [R[f] - R_{\text{emp}}[f]] \leq \epsilon_{\mathcal{F}} \right] \geq 1 - \delta$$

1. Fix δ
2. Compute $f_{\mathcal{F}}^* = \arg \min_{f \in \mathcal{F}} [R_{\text{emp}}[f] + \epsilon_{\mathcal{F}}]$ for a sequence of spaces
 $\mathcal{F}_1 \subset \mathcal{F}_2 \subset \dots \subset \mathcal{F}_k$
3. Return f_i^* with smallest $R_{\text{emp}}[f_i^*] + \epsilon_{\mathcal{F}_i}$

Does this work?

The problem with error bounds

Most bounds are hopelessly loose. Typically, we get for $1 - \delta = .95$

$$\epsilon = 3000.$$

Main culprit is the uniformity requirement. Can we still use them for anything or are they just a weird sport?

Form of bounds is important, even if their value is not. In particular, **large margin** is good.

Hilbert Space Methods

SVM's: the old story

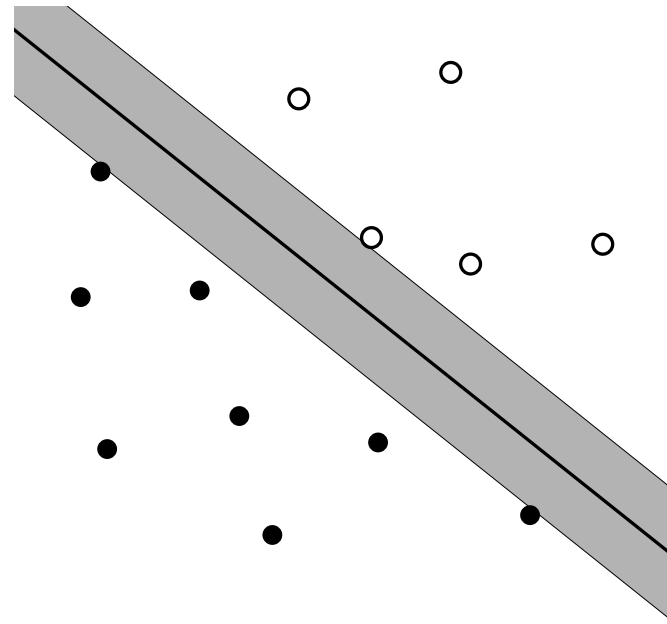
Kernel $k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ pos.def.
similarity measure

Feature map $\Phi : \mathcal{X} \mapsto \mathcal{F}$
obeys $k(x, x') = \langle \Phi(x), \Phi(x') \rangle$

e.g. Gaussian Kernel:

$$k(x, x') = e^{-\|x-x'\|^2/(2\sigma^2)}$$

Find maximum margin separating hyper-
plane in high dimensional space!



$$f(x) = \text{sgn} \left[b + \sum_{i=1}^m \alpha_i k(x_i, x) \right]$$

Want more general story behind Hilbert space methods. How do we tell what is a good kernel, anyway? Want large margin. What kernel will give us large margin?

Lessons so far:

- capacity control is crucial;
- large margin is good;
- pursue abstract approach looking for general $f : \mathcal{X} \mapsto \mathcal{Y}$, worry about actual algorithm later.

Regularized Risk

Motivated by form of error bounds, minimize

$$R_{\text{reg}}[f] = \underbrace{\frac{1}{m} \sum_{i=1}^m L(f(x_i), y_i)}_{R_{\text{emp}}[f]} + \underbrace{\Omega[f]}_{\text{regularizer}}$$

over some large space of functions \mathcal{H} .

$\Omega[f]$ is a penalty term penalizing hypotheses that are too complex. Effectively SRM.

See Regularization networks of Poggio & Girosi.

Regularized Spaces of Functions

Given $\{(x_1, y_1), \dots, (x_m, y_m)\}$ look for $f : \mathcal{X} \mapsto \mathcal{Y}$ in some linear space of functions \mathcal{H} minimizing

$$R_{\text{reg}}[f] = \frac{1}{m} \sum_{i=1}^m L(f(x_i), y_i) + \|f\|_{\mathcal{H}}^2$$

Regularized Spaces of Functions

Given $\{(x_1, y_1), \dots, (x_m, y_m)\}$ look for $f : \mathcal{X} \mapsto \mathcal{Y}$ in some linear space of functions \mathcal{H} minimizing

$$R_{\text{reg}}[f] = \frac{1}{m} \sum_{i=1}^m L(f(x_i), y_i) + \|f\|_{\mathcal{H}}^2$$

$$R_{\text{reg}}[f] = \frac{1}{m} \sum_{i=1}^m L(f(x_i), y_i) + \langle f, f \rangle_{\mathcal{H}} \quad \text{Hilbert space}$$

Regularized Spaces of Functions

Given $\{(x_1, y_1), \dots, (x_m, y_m)\}$ look for $f : \mathcal{X} \mapsto \mathcal{Y}$ in some linear space of functions \mathcal{H} minimizing

$$R_{\text{reg}}[f] = \frac{1}{m} \sum_{i=1}^m L(f(x_i), y_i) + \|f\|_{\mathcal{H}}^2$$

$$R_{\text{reg}}[f] = \frac{1}{m} \sum_{i=1}^m L(f(x_i), y_i) + \langle f, f \rangle_{\mathcal{H}}$$

Hilbert space

$$R_{\text{reg}}[f] = \frac{1}{m} \sum_{i=1}^m L(\langle f, k_{x_i} \rangle, y_i) + \langle f, f \rangle_{\mathcal{H}}$$

RKHS

Regularized Spaces of Functions

Given $\{(x_1, y_1), \dots, (x_m, y_m)\}$ look for $f : \mathcal{X} \mapsto \mathcal{Y}$ in some linear space of functions \mathcal{H} minimizing $R_{\text{reg}}[f]$.

$$R_{\text{reg}}[f] = \frac{1}{m} \sum_{i=1}^m L(f(x_i), y_i) + \|f\|_{\mathcal{H}}^2$$

$$R_{\text{reg}}[f] = \frac{1}{m} \sum_{i=1}^m L(f(x_i), y_i) + \langle f, f \rangle_{\mathcal{H}} \quad \text{Hilbert space}$$

$$R_{\text{reg}}[f] = \frac{1}{m} \sum_{i=1}^m L(\langle f, k_{x_i} \rangle, y_i) + \langle f, f \rangle_{\mathcal{H}} \quad \text{RKHS}$$

The k_x are “prototypical” functions s.t. $f(x) = \langle f, k_x \rangle$.

Representer Theorem

Minimizer of

$$R_{\text{reg}}[f] = \frac{1}{m} \sum_{i=1}^m L(\langle f, k_{x_i} \rangle, y_i) + \langle f, f \rangle_{\mathcal{H}}$$

will be in the span of $k_{x_1}, k_{x_2}, \dots, k_{x_m}$!

The hypothesis can be written

$$f(x) = \langle f, k_x \rangle = \sum_{i=1}^m \alpha_i \langle k_{x_i}, k_x \rangle = \sum_{i=1}^m \alpha_i k(x, x_i).$$

where $k(x, x') = \langle k_x, k_{x'} \rangle$.

All we need to find are $\alpha_1, \alpha_2, \dots, \alpha_m$. How do we construct the RKHS?

Constructing the RKHS

$$f(x) = \langle f, k_x \rangle$$

Bootstrap everything from $k(x, x') = \langle k_x, k_{x'} \rangle$ for $x, x' \in \mathcal{X}$!

1. Anything outside $\text{span} \{ k_x \mid x \in \mathcal{X} \}$ is uninteresting, so $f = \int_{\mathcal{X}} \beta(x) k_x dx$.

2. To evaluate $f(x')$ use

$$f(x) = \langle f, k_{x'} \rangle = \int_{\mathcal{X}} \beta(x) \langle k_x, k_{x'} \rangle dx = \int_{\mathcal{X}} \beta(x) k(x, x') dx$$

3. To compute $\langle f, f' \rangle$ use

$$\langle f, f' \rangle = \int_{\mathcal{X}} \int_{\mathcal{X}} \beta(x) \beta(x') \langle k_x, k_{x'} \rangle dx dx' = \int_{\mathcal{X}} \int_{\mathcal{X}} \beta(x) \beta(x') k(x, x') dx dx'$$

4. Note that $k_x(x') = \langle k_x, k_{x'} \rangle = k(x, x')$ so we simply have $k_x = k(x, x')$.

5. \mathcal{H} is a particular instance of a feature space \mathcal{F} if we set $\Phi(x) = k_x$.

Correspondence

$$R_{\text{reg}}[f] = \frac{1}{m} \sum_{i=1}^m L(\langle f, k_{x_i} \rangle, y_i) + \langle f, f \rangle_{\mathcal{H}}$$

Kernel methods make sense from Regularization Theory point of view if kernel corresponds to sensible operator $\Omega[f] = \langle f, f \rangle_{\mathcal{H}}$.

Fourier regularization

Fourier transform on \mathbb{R}^n :

$$\hat{f}(\omega) = \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} f(x) e^{-i\omega \cdot x} dx$$

Inverse transform:

$$f(x) = \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{C}^n} \hat{f}(\omega) e^{i\omega \cdot x} d\omega$$

Fourier regularization:

$$\Omega[f] = \langle f, f \rangle_{\mathcal{H}} = \int e^{\sigma^2 \|\omega\|^2 / 2} \hat{f}(\omega)^2 d\omega$$

Corresponding kernel:

$$k(x, x') = e^{-\|x - x'\|^2 / (2\sigma^2)}$$

The Gaussian kernel will heavily penalize non-smooth functions!

Other kernels

- Homogeneous polynomial: $k(x, x') = (x \cdot x')^p$
- Non-homogeneous polynomial: $k(x, x') = (x \cdot x' + 1)^p$
- tanh kernel: $k(x, x') = \tanh(\kappa(x \cdot x') + \delta)$
- Triangular kernel: $k(x, x') = 1 - |(x - x') / d|$
- String kernels: $k(\text{string}_1, \text{string}_2)$
- Kernels on distributions: Fisher, etc.
- Diffusion kernels: $k(x, x') = [e^{\beta\Delta}]_{x, x'}$

Similarity measure \longleftrightarrow Regularization

Algorithms

Modularity of Hilbert space methods

$$f^* = \arg \min_{f \in \mathcal{H}} \left[\frac{1}{m} \sum_{i=1}^m \underbrace{L(\langle f, k_{x_i} \rangle, y_i)}_{\text{Determines algorithm}} + \underbrace{\langle f, f \rangle_{\mathcal{H}}}_{\text{Determines kernel}} \right]$$

- Same algorithm (SVM) can be used in very different contexts by changing the kernel \rightarrow kernel engineering
- Regularization scheme can be studied independent of application (classification, regression, etc.)
- ANY kernel method can be formulated as one of these minimization problems

Soft margin SVM's

Relax problem to learning continuous functions $f : \mathcal{X} \mapsto \mathbb{R}$ with hinge loss

$$L(f(x), y) = C \max(0, -y f(x) + 1)$$

Then

$$f^* = \arg \min_{f \in \mathcal{H}} \left[\frac{1}{m} \sum_{i=1}^m L(\langle f, k_{x_i} \rangle, y_i) + \langle f, f \rangle_{\mathcal{H}} \right]$$

reduces to soft margin SVM

$$f^* = \arg \min_{f \in \mathcal{H}} \left[\langle f, f \rangle + C \sum_{i=1}^m \xi_i \right] \quad \text{subject to} \quad y_i f(x_i) \geq 1 - \xi_i$$

Probably the most popular algorithm for classification.

Kernel Regression

If we set then

$$f^* = \arg \min_{f \in \mathcal{H}} \left[\frac{1}{m} \sum_{i=1}^m L(\langle f, k_{x_i} \rangle, y_i) + \langle f, f \rangle_{\mathcal{H}} \right]$$

reduces to soft kernel regression

$$f^* = \arg \min_{f \in \mathcal{H}} \left[\langle f, f \rangle + C \sum_{i=1}^m (\xi_i + \xi'_i) \right] \quad \text{subject to} \quad \begin{aligned} y_i - f(x_i) &\leq \epsilon + \xi \\ y_i - f(x_i) &\geq -\epsilon - \xi' \end{aligned}$$

