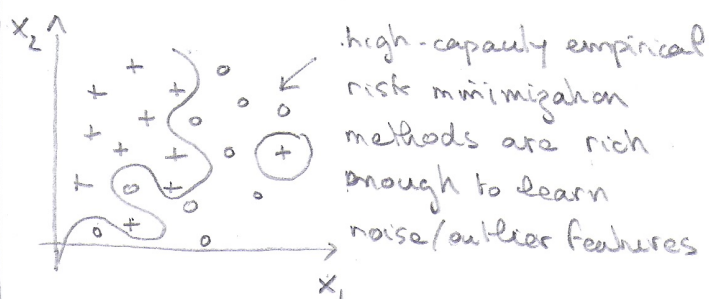
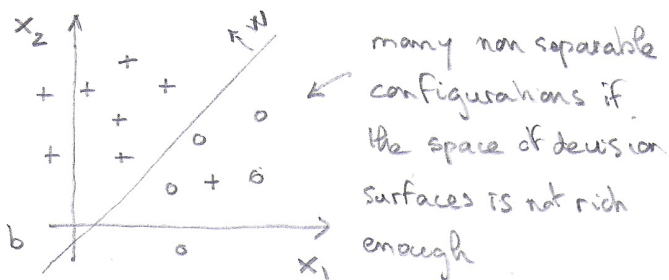


# Motivation For SVM

①

\* Two main limitations to statistical learning algorithm (before GPU):



## Capacity limitation

↳ idea for solution:

high-dimensional embedding  
for a given number of samples configurations become separable in high dimension (assuming general position)

## Overfitting limitation

↳ idea for solution:

Tikhonov regularization  
taking irregularity of decision surfaces to be the signature of overfitting, we restrain the space of candidate surfaces by penalizing irregularity.

\* Claim: Support vector machines (SVMs) resolve the above tradeoff by implementing a computationally efficient high-dimensional embedding that has good regularity property.

\* Theoretical Foundation: Reproducing kernel Hilbert space (RKHS)

Main idea: choosing a positive definite kernel function  $k: X \times X \rightarrow \mathbb{R}$  defined on the data space  $X$  also

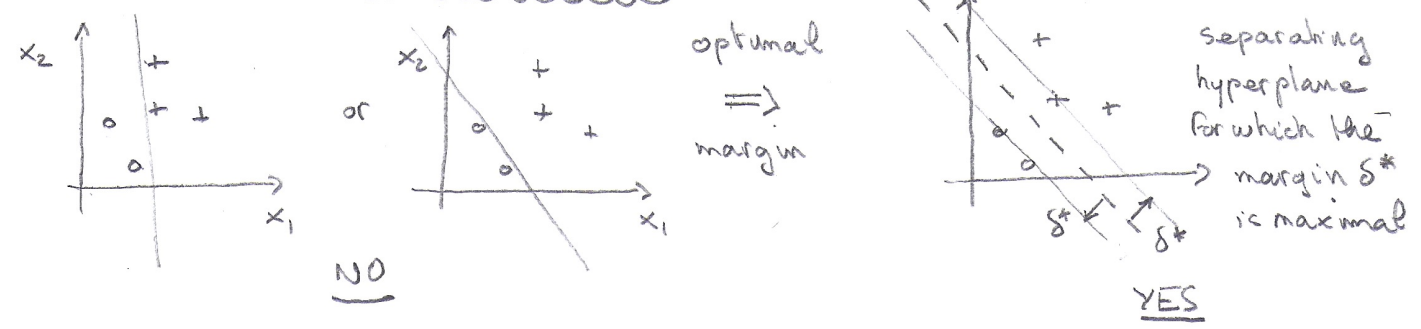
defines a high-dimensional (infinite dimensional) embedding in a RKHS  $H$  via  $x \mapsto (\forall y \mapsto k(x, y))$ .

Computationally: The kernel function  $k$  measures similarity between data point once embedded and behaves as a dot product:  $k(x, y) = \langle k(x, \cdot), k(y, \cdot) \rangle_H \leftarrow$  Hilbert dot product

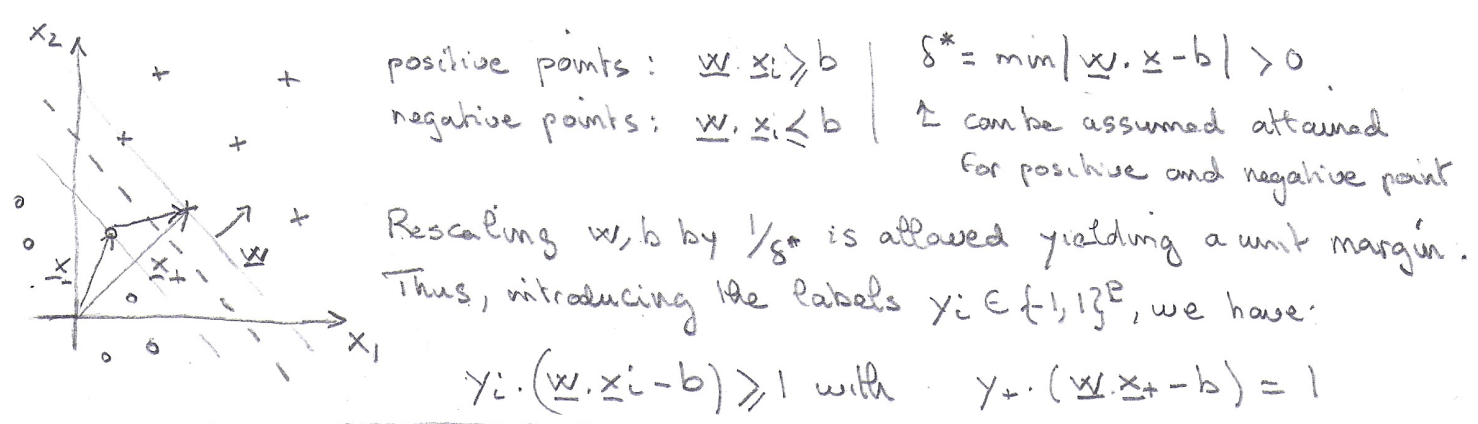
# Optimal margin classifier

- \* Algorithms that can exploit the kernel trick need to be formulated in terms of dot products on sample points  $\underline{x}_i, i=1, \dots, P$ .
- \* Optimal margin classifiers provide such a formulation for binary linear classification (a.k.a perceptron)

What is the "best" classification?



Constrained optimization problem:



$$y_i \cdot (\underline{w} \cdot \underline{x}_i - b) \geq 1 \quad \text{with} \quad y_+ \cdot (\underline{w} \cdot \underline{x}_+ - b) = 1$$

$$y_i \cdot (\underline{w} \cdot \underline{x}_- - b) = 1$$

Margin magnitude

$$\frac{\underline{w}}{\|\underline{w}\|} \cdot (\underline{x}_+ - \underline{x}_-) = \frac{b+1 - (b-1)}{\|\underline{w}\|} = \frac{2}{\|\underline{w}\|}$$

If  $\underline{x}_+, \underline{x}_-$  are on the positive and negative boundaries respectively.

Optimal margin:

$$\underline{w}^* = \arg \max_{\underline{w}} \frac{2}{\|\underline{w}\|} \quad \text{such that} \quad y_i \cdot (\underline{w} \cdot \underline{x}_i) \geq 1$$

$$= \arg \min_{\underline{w}} \frac{\|\underline{w}\|^2}{2} \quad \text{such that} \quad y_i \cdot (\underline{w} \cdot \underline{x}_i) \geq 1$$

This is a quadratic optimization problem subjected to linear constraints.

## Lagrangian Formulation

3

To include the constraints in the optimization of the margin we consider the following Lagrangian function

$$L(\underline{w}, \underline{\alpha}) = \frac{1}{2} \|\underline{w}\|^2 - \sum_i \alpha_i \left[ y_i (\underbrace{\underline{w} \cdot \underline{x}_i}_{\text{def } w_0} - b) - 1 \right], \quad \alpha_i \geq 0$$

↑ Lagrange multiplier

$$\text{We have: } \max_{\underline{\alpha} \geq 0} L(\underline{w}, \underline{\alpha}) = \begin{cases} +\infty & \text{if } y_i (\underline{w} \cdot \underline{x}_i - b) < 1 \text{ for some } i \\ \frac{1}{2} \|\underline{w}\|^2 & \text{otherwise} \end{cases}$$

$$\text{Thus we seek to find: } \underline{w}^* = \arg \left( \min_{\underline{w}} \max_{\underline{\alpha} \geq 0} L(\underline{w}, \underline{\alpha}) \right)$$

=  $\underline{I}^*$

This is the primal problem. Let us consider the dual problem obtained by interverting min and max:

$$\underline{L}^* = \max_{\underline{\alpha} \geq 0} \min_{\underline{w}} L(\underline{w}, \underline{\alpha})$$

In general  $\underline{L}^* \leq \underline{I}^*$ . Indeed, posing

$$\begin{cases} \underline{L}^* = \min_{\underline{w}} L(\underline{w}, \underline{\alpha}^*) \\ \underline{I}^* = \max_{\underline{\alpha}} L(\underline{w}^*, \underline{\alpha}) \end{cases}$$

$$\text{we have: } \underline{L}^* = \min_{\underline{w}} L(\underline{w}, \underline{\alpha}^*) \leq L(\underline{w}^*, \underline{\alpha}^*) \leq \max_{\underline{\alpha}} L(\underline{w}^*, \underline{\alpha}) = \underline{I}^*$$

The duality gap  $\underline{I}^* - \underline{L}^*$  is zero under certain convexity conditions that are met in our situation.

Thus solving the dual problem (easier) is equivalent to solving the primal problem.

# SVM method

Writing the KKT conditions for the dual problem yields

$$\frac{\partial L}{\partial \underline{w}} = 0 \Rightarrow * \underline{w} - \sum_i \alpha_i \gamma_i \underline{x}_i = 0 \Leftrightarrow \underline{w} = \sum_i \alpha_i \gamma_i \underline{x}_i$$

$$* \sum_i \alpha_i \gamma_i = 0$$

↑  
The optimal weight vector is a linear combination of some positions vectors ( $\alpha_i \neq 0$ ): the support vectors.

The dual problem becomes finding the  $\alpha_i$ 's, i.e., solving

$$\arg \max_{\alpha \geq 0} \left\{ \frac{1}{2} \left( \sum_i \alpha_i \gamma_i \underline{x}_i \right) \left( \sum_j \alpha_j \gamma_j \underline{x}_j \right) - \sum_i \alpha_i \left[ \gamma_i \left( \sum_j \alpha_j \gamma_j \underline{x}_j \right) \underline{x}_i - \gamma \right] \right\}$$
  
$$\sum_i \alpha_i \gamma_i = 0$$

$$= \arg \max_{\alpha \geq 0} \left\{ -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j \gamma_i \gamma_j (\underline{x}_i \cdot \underline{x}_j) \right\}$$
  
$$\sum_i \alpha_i \gamma_i = 0$$

↑ dot product on the data vectors  $\Rightarrow$  kernel trick possible!

$\underline{x}_i \cdot \underline{x}_j \leftrightarrow k(\underline{x}_i, \underline{x}_j)$  where  $k$  is some well-chosen positive definite function

The problem becomes a problem for numerical analysts

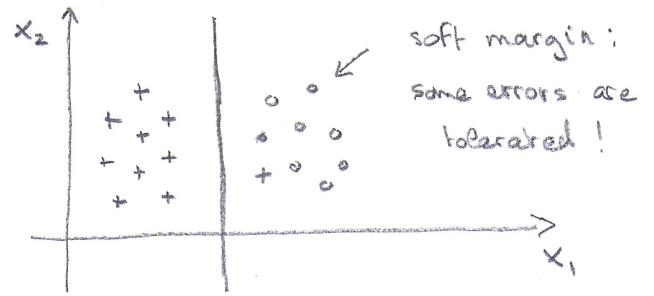
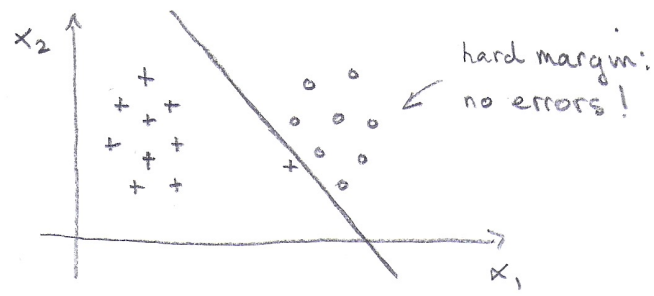
The existence and uniqueness of a solution is guaranteed by convexity under assumption of separability.

Popular kernel choice:  $k(x, y) = (1 + xy)^d \leftarrow$  polynomials of degree up to  $d$ .

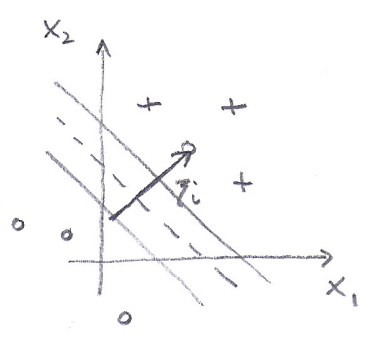
$k(x, y) = e^{-\frac{\|x-y\|}{\sigma}} \leftarrow$  radial basis function

if  $\sigma \ll 1 \rightarrow$  tend to overfit. tuning is involved.

Link to regularization



Idea: assigning a cost to errors via new optimization variables called slackness variables  $\zeta_i, i=1, \dots, P$ .



$$\zeta_i = \left[ 1 - \gamma_i (\underline{w} \cdot \underline{x}_i - b) \right]_+ \leftarrow \text{positive part}$$

thus  $\zeta_i \geq 0 \leftarrow \text{constraints}$

Thus suggests assigning the cost  $\sum_i \zeta_i$

Soft margin optimization:

trade off tuning variable

Primal problem:  $\arg \min_{\underline{w}, \zeta \geq 0} \left( \frac{1}{2} \|\underline{w}\|^2 + C \sum_i \zeta_i \right)$

such that:  $\gamma_i (\underline{w} \cdot \underline{x}_i - b) \geq 1 - \zeta_i$

Dual problem:  $\arg \min_{0 \leq \alpha_i \leq C} \left( \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j \gamma_i \gamma_j (\underline{x}_i \cdot \underline{x}_j) - \sum_i \alpha_i \right)$   
 $\sum_i \alpha_i \gamma_i = 0$

Formulation in term of regularization

$$w^* = \arg \min_w \left( \underbrace{C(w, x_i, y_i)}_{\text{empirical risk}} + \frac{1}{2C} \underbrace{\|w\|^2}_{\text{norm on the separating vectors}} \right)$$

$$F^* = \arg \min_{F \in H} \left[ C(F, x_i, y_i) + \lambda \Omega(\|F\|_H^2) \right]$$

$\uparrow$  RKHS space obtained via the choice of a kernel  
 $\uparrow$  tradeoff parameter  
 $\uparrow$  RKHS norm penalizing irregular decision surfaces  
 $\uparrow$  strictly increasing function

## Representer Theorem

A priori, regularization problem are infinite dimensional.

The representer theorem state that it is not the case in RKHS.

$$\arg \min_{F \in H} [C(F, x_i, y_i) + \lambda \Omega(\|F\|_H^2)] \in \left\{ \sum_{i=1}^P \alpha_i K(x_i, \cdot) \mid \alpha_i \in \mathbb{R} \right\} = F$$

↑ finite dimensional.

### Sketch of the proof

For  $F$  in  $H$ ,  $F = F_{\parallel} + F_{\perp}$ ,  $F_{\parallel} \in F$  and  $F_{\perp} \in F^{\perp}$ : orthogonal space to the finite dimensional space  $F$ .

For all  $i$ ,  $\langle F_{\perp}, K(x_i, \cdot) \rangle = F_{\perp}(x_i) = 0$ .

Thus, for all  $F_{\perp}$ ,  $C(F, x_i, y_i) = C(F_{\parallel}, x_i, y_i)$  while

$$\|F\|_H^2 = \|F_{\parallel}\|_H^2 + \|F_{\perp}\|_H^2 \gg \|F_{\parallel}\|_H^2.$$

## Regularization via Mercer's Theorem

Consider a definite positive kernel  $k: X \times X \rightarrow \mathbb{R}$  such that

$\int_X \int_X k(y, x)^2 < +\infty$  (For all  $F, g$  in  $L^2(X)$ , we have

$\int_X \int_X k(y, x) F(y) F(x) \geq 0$ .) This defines an operator

$$T_k: L^2(X) \rightarrow L^2(X), T_k[F](y) = \int_X k(y, x) F(x) dx.$$

Then:  $T_k$  has a spectral decomposition: there is a complete system of orthonormal functions  $e_i, i \in \mathbb{N}$ .

such that  $K(y, x) = \sum_i \lambda_i e_i(y) e_i(x)$ ,  $\lambda_i \rightarrow 0$   
 ↑ eigenvalues  $\sum_i \lambda_i^2 < +\infty$

The above theorem allows one to specify a high-dimensional embedding into the RKHS associated to  $K$ .

Consequences

Posit:  $\phi(x) = \sum_i \sqrt{\lambda_i} e_i(x) e_i(\cdot) \in H \subset L_2(X)$

We have:  $\langle \phi(x), \phi(y) \rangle_{L_2(X)} = \sum_i \sum_j \sqrt{\lambda_i \lambda_j} e_i(x) e_j(y) \underbrace{\langle e_i, e_j \rangle}_{\delta_{ij}}$   
 natural dot product  $\uparrow$   $= \sum_i \lambda_i e_i(x) e_i(y)$   
 $= K(x, y) \leftarrow$  reproducing kernel

How to relate  $\langle, \rangle_H$  to  $\langle, \rangle_{L_2}$ ?

By finding  $\langle, \rangle_H$  such that  $K(x, y) = \langle K(x, \cdot), K(y, \cdot) \rangle_H$  holds.

Posit  $\langle f, g \rangle_H = \sum_i \frac{1}{\lambda_i} \int_x e_i(x) f(x) dx \int_x e_i(x) g(x) dx$

$\langle K(x, \cdot), K(y, \cdot) \rangle_H = \sum_i \frac{1}{\lambda_i} \lambda_i e_i(x) \lambda_i e_i(y) = K(x, y)$

$f \in H, \|f\|_H^2 = \sum_i \frac{(\int e_i(x) f(x) dx)^2}{\lambda_i}$   $\leftarrow$  component of  $f$  on  $e_i$  coefficients must decay to 0 faster than  $\lambda_i$  as  $i \rightarrow +\infty$

High-frequency components are penalized by  $\| \cdot \|_H$

Translation-invariant kernel

$T_x[f](y) = \int_x k(y-x) f(x) dx$ , eigenfunction  $f_\omega(x) = e^{i\omega x}$

$T_x[f_\omega](y) = e^{i\omega y} \int_x k(y-x) e^{i\omega(x-y)} dx = \hat{K}(\omega) f_\omega(y)$   
 $\uparrow$  Fourier transform

$\phi(x) = \left( \int_0^{+\infty} \sqrt{\hat{K}(\omega)} e^{i\omega x} d\omega \right) e^{i\omega(\cdot)}$   $\leftarrow$  argument  
 $\uparrow$  positive measure by Bochner theorem

$\langle f, g \rangle_H = \int_0^{+\infty} d\omega \frac{\overline{\hat{f}(\omega)} \hat{g}(\omega)}{\hat{K}(\omega)}$ ,  $\|f\|_H^2 = \int_0^{+\infty} d\omega \frac{\|\hat{f}(\omega)\|^2}{\hat{K}(\omega)}$   
 $\uparrow$  penalize irregularity  $\omega \rightarrow +\infty$