# Non-linear discriminant functions.

$$\phi : R^m \rightarrow R^{m'}$$

New training set,

$$\{(z_i, y_i), i = 1, \ldots, n\}, \quad z_i = \phi(x_i).$$

New Dual:

$$\max_{y} \; q(u) = \sum_{i=1}^{n} u_i - \frac{1}{2} \sum_{i,j=1}^{n} u_i u_j \, y_i y_j \, \phi(x_i)^T \phi(x_j)$$

$$\text{s.t} \quad 0 \leq u_i \leq C, \quad i = 1, \ldots, n, \quad \sum_{i=1}^{n} y_i u_i = 0$$

The problem is still a QP problem over $R^n$ irrespective of $\phi$ & $m'$.

But we still want to compute $\phi(x)$.

(187)

# Kernel idea.

Suppose $\exists$ a fn. $K: R^m \times R^m \rightarrow R$ s.t

$$K(x_i, x_j) = \phi(x_i)\phi(x_j)$$

& computing $K(x_i, x_j)$ is as expensive as $x_i^T x_j$

then, Dual can be solved by replacing

$z_i^T z_j$ by $K(x_i, x_j)$.

What happens during testing?

we have, $W^* = \sum \alpha_i^* y_i \phi(x_i)$

& $b^* = y_j - \phi(x_j)^T w^* = y_j - \sum_i \alpha_i^* y_i \phi(x_i)^T\phi(x_j)$

$\forall$ test pattern $x$, we need to compute.

$$f(x) = \emptyset \phi(x)^T w^* + b.$$

$$= \sum_i \alpha_i^* y_i \phi(x_i)^T \phi(x) + b^*$$

$$= \sum_i \alpha_i^* y_i K(x_i, x) + \left(y_j - \sum_i \alpha_i^* y_i K(x_i, x_j)\right)$$

(188)

∴ Never do we need to compute $\phi$

In theory, $\phi$ can even by $\infty$ dim.

∴ For an SVM, all that is needed to be

stored is

$\mu_i^*$ & $x_i^*$ for $i \in S$.

       ↓

      Support
      vectors.

---

Eg of a Kernel fn in $R^2$ : $k(x_i, x_j)$

Let $x \in R^2$, $x_i = (x_{i1}, x_{i2})^T$     $= (1 + x_i^T x_j)^2$

$$k(x_i, x_j) = (1 + x_{i1}x_{j1} + x_{i2}x_{j2})^2$$

To show, $\exists \phi$ in $m' > m$    s.t $\phi(x)^T \phi(x)$

                                 $= k(x_i, x_j)$.

Consider $\phi : R^2 \to R^6$

$$\phi(x) = \begin{bmatrix} 1 & \sqrt{2}x_1 & \sqrt{2}x_2 & x_1^2 & x_2^2 & \sqrt{2}x_1 x_2 \end{bmatrix}$$

One can show that $\phi(x_i)^T \phi(x_j) = k(x_i, x_j)$

Note: $\phi$ is non-unique.

(189...

# Kernels in general.

Mercer theorem: Given a symmetric fn

$$K: \mathbb{R}^m \times \mathbb{R}^m \to \mathbb{R}, \quad \exists \text{ an inner product}$$

space $H$ & mapping $\phi: \mathbb{R}^m \to H$, so that

$$K(x_1, x_2) = \phi(x_1)^T \phi(x_2) \text{ if for all sq. integrable}$$

fns $g$,

$$\int K(x_1, x_2) \, g(x_1) g(x_2) \, dx_1 dx_2 \geq 0$$

In other words, if $\bar{K}$ a $n \times n$ matrix

with $\bar{K}_{ij} = K(x_i, x_j)$. If $\bar{K}_{nn}$ is PSD for all

$n$ data points, then $K$ is a valid kernel.
(set)

$$\sum_{i,j=1}^{n} c_i c_j \, K(x_i, x_j) \geq 0$$

(190)

One can show that

1) Polynomial kernel:

$$K_P(x_1, x_2) = \left(1 + x_1^T x_2\right)^P$$

2) Gaussian kernel:

$$K_G(x_1, x_2) = e^{-\frac{\|x_1 - x_2\|^2}{\sigma^2}}$$

3) sigmoidal kernel:

$$K_S(x_1, x_2) = \tanh\left(a x_1^T x_2 + \theta\right)$$

all the above satisfy Mercer's theorem.

SVM with $K_G$:

$$f(x) = \sum_{i \in S} \hat{y}_i y_i K(x_i, x) + b^*$$

$$= \sum_{i \in S} \hat{y}_i y_i e^{-\frac{\|x_i - x\|^2}{2\sigma^2}} + b^*$$

[RBF NN]

(191).

## SVM with Ks

$$f(x) = \sum_{i \in S} \alpha_i^* y_i \tanh\left(a x^T x_i + \theta\right) + b^*$$

[NN with one hidden layer with tanh activation
# of nodes in hidden determined by $\alpha_i^*$]

Why do SVMs perform well:

$$E\,P_{err} \leq \min\left(\frac{S}{n}, \frac{R^2 \|w\|^2}{n}, \frac{m}{n}\right).$$

$S \rightarrow$ no of support vectors

$R \rightarrow$ radius of smallest sphere enclosing all example.

$\|w\|^{-2} \rightarrow$ margin of the hyper-plane

$m \rightarrow$ feature dim

$n \rightarrow$ no of examples.

1) Good data compression

2) large margin

3) dim of feature space is small.

(192).

## SVM from a risk min view.

we have
$$\min_{w, b, \xi} \frac{1}{2} w^T w + c \sum_{i=1}^{n} \xi_i$$

$$\text{s.t} \quad y_i(w^T x_i + b) \geq 1 - \xi_i, \quad i = 1, \cdots, n$$

$$\xi_i \geq 0, \quad i = 1, \cdots, n$$

Now given any $w, b$, $\xi_i$ has to satisfy the following.

$$\xi_i \geq \max\left(0, 1 - y_i(w^T x_i + b)\right)$$

∴ The above problem can be effectively written as

$$\min_{w, b} \frac{1}{2} w^T w + c \sum_{i=1}^{n} \max\left(0, 1 - y_i(w^T x_i + b)\right)$$

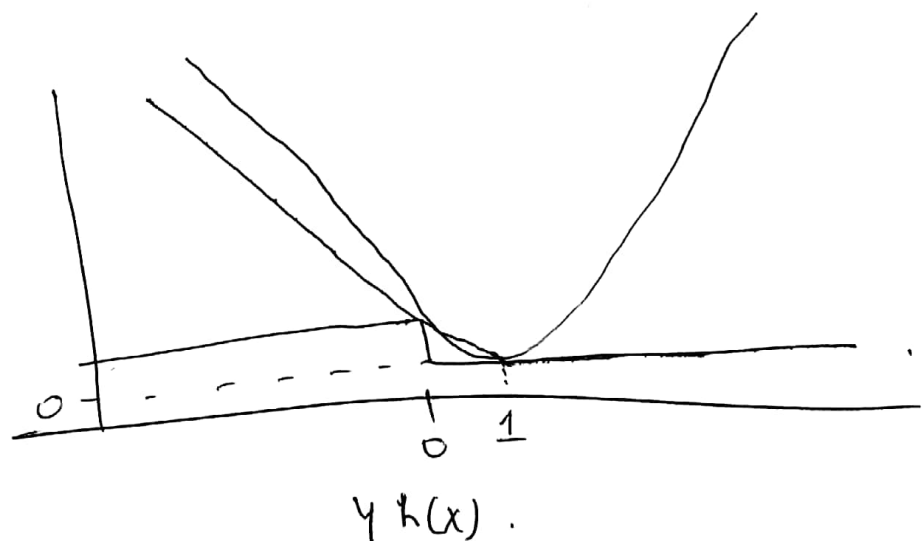& final classifier, $f(x) = w^T x + b$.

We know that 0-1 loss is non-differentiable.

$L(h(x), y)$ can be made into a fn of single variable $yf(x)$, if $y \in \{-1, 1\}$.

For 0-1 loss, $yh(x)$ is

$$L_{01} = 1 \quad \text{if} \quad yh(x) < 0$$
$$\quad\quad\quad 0 \quad \text{otherwise}$$

$$L_{sq\text{-error}} = \left(1 - yh(x)\right)^2$$

$$L_{hinge} = \max\left(0, 1 - yh(x)\right)$$



$yh(x)$.

∴ under this formulation all losses are Convex approxim. for 0-1 loss.

(1a) .

∴ SVM optimization problem can be written as follows:

$$\min_{w, b} \frac{1}{n} \sum_{i=1}^{n} L\left(y_i, f(x_i)\right) + c' \frac{1}{2} w^T w$$

$$f(x_i) = w^T x_i + b.$$

∴ SVM is empirical risk minimization under hinge-loss with $L_2$ regularizer.

[soft-margin loss].

How would kernels fit in this framework?

(195).

# Representer theorem:

For any positive definite kernel, $\exists$ a vector space with an inner product $\mathcal{H}$, s.t kernel is the inner product in that space.

Mercer thorem says that $\exists$ $\phi$ from $X$ to $\mathcal{H}$.

with these, representer theorem says

Let $\Omega : [0, \infty) \to R^+$ be a strictly monotonically increasing function. Then any minimizer $g$ over $\mathcal{H}$ of the regularized risk

$$ C\Big(\big(x_i, y_i, g(x_i)\big), i = 1, \ldots n\Big) + \Omega\left(\|g\|^2\right). $$

admits a representation

$$ g(x) = \sum_{i=1}^{n} \alpha_i \, k(x_i, x). $$

$$ (196) $$

This is a very powerful theorem because
it says that the minimizer of the
empirical risk $\hat{G}$ is a linear combination
of kernels centered around data points alone‼

∴ Even though $\mathcal{H}$ may be very high dim.
One can design an opti problem for risk
minimization by searching for $n$ real no's $\alpha_i$

This is preciselly what SVM does‼

---

The idea of kernel is generalizable.

Eg: suppose we are so doing Knn in $\phi(x)$
dim.

Let $C_+ = \dfrac{1}{n_+} \sum_{i: y_i = +1} \phi(x_i)$ , $C_- = \dfrac{1}{n_-} \sum_{i: y_i = -1} \phi(x_i)$

(197).

Now, KNN would put a new pattern in class
+1, if

$$\| \phi(x) - c_+ \|^2 < \| \phi(x) - c_- \|^2$$

$$\| \phi(x) - c_+ \|^2 = \phi(x)^T \phi(x) - 2\phi(x)^T c_+ + c_+^T c_+$$

$\Rightarrow$ we put $x$ in class +1 if

$$\bullet \quad \phi(x)^T c_+ - \phi(x)^T c_- + \frac{1}{2}\left(c_-^T c_- - c_+^T c_+\right) > 0$$

$$\phi(x)^T c_+ = \phi(x^T) \left(\frac{1}{n_+} \sum_{i: y_i = +1} \phi(x_i)\right)$$

$$= \frac{1}{n_+} \sum K(x_i, x)$$

$$c_+^T c_+ = \frac{1}{n_+^2} \sum_{i,j: y_i = y_j = 1} K(x_i, x_j)$$

very much related to kernel density estimates.