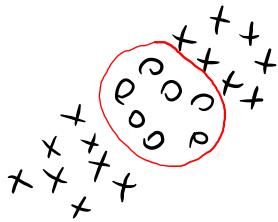


SUPPORT VECTOR MACHINES

27/4/21



Main Idea

- ① Can we augment the feature space so that the data become linearly separable in the higher dimensional space?
- ② Can we do this without having to work on the higher dimensional space?

MOTIVATION

p features : x_1, \dots, x_p

Augmented 2p features : $x_1, \frac{x_1^2}{y_1}, x_2, \frac{x_2^2}{y_2}, \dots, x_p, \frac{x_p^2}{y_{2p}}$

$x_1 x_2$

$x_i x_k$

$x_1 x_2 x_3$

SUPPORT VECTOR
CLASSIFIER ON

$$\min_{\beta_0, \beta_1, \xi} \frac{1}{2} \|\beta\|^2$$

$$\text{s.t. } y_i (\beta_0 + \sum_{j=1}^p \beta_{j1} x_{ij} + \sum_{j=1}^p \beta_{j2} x_{ij}^2) \geq 1 - \xi_i \quad \forall i$$

$$\xi_i \geq 0$$

$$\sum_{i=1}^n \xi_i \leq C$$

In the augmented space : LINEAR DECISION SURFACE $\alpha_0 + \alpha^T y = 0$

In the original space : $\beta_0 + \sum_{j=1}^p \beta_{j1} x_{ij} + \sum_{j=1}^p \beta_{j2} x_{ij}^2 = 0$ QUADRATIC
DECISION SURFACE

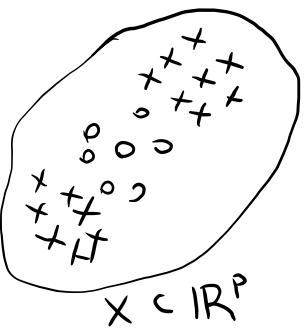
SVC : Linear decision surface

$$f(\underline{x}) = \beta_0 + \underline{\beta}^t \underline{x} = \beta_0 + \sum_{j=1}^p \beta_j x_j = \beta_0 + \underline{\beta} \cdot \underline{x}$$
$$\leftarrow = \beta_0 + \sum_{i=1}^n \gamma_i y_i \underline{x}_i \cdot \underline{x}$$
$$\underline{\beta} = \sum_{i=1}^n \gamma_i y_i \underline{x}_i$$

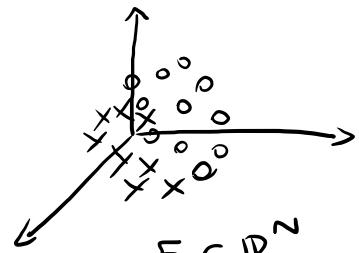
Also in the estimation of parameters, only inner products appear:

$$\max_{\underline{\gamma}} \sum_{i=1}^n \gamma_i - \frac{1}{2} \sum_{i,j=1}^n \gamma_i \gamma_j y_i y_j \underline{x}_i \cdot \underline{x}_j$$

(2) inner products



ϕ



$$f(x) = \beta_0 + \beta^t \phi(x)$$

$$\beta = \sum_{i=1}^n \gamma_i y_i \phi(x_i)$$

$$f(x) = \beta_0 + \sum_{i=1}^n \gamma_i y_i \phi(x_i) \cdot \phi(x)$$

For estimation, we would need

$$\phi(x_i) \cdot \phi(x_j)$$

KERNEL TRICK

Find a function (Kernel) which is defined in the input space but which acts as a dot product in some feature space \mathbb{R}^N for some map ϕ .

$$\underline{K(x, z) = \phi(x) \cdot \phi(z)} \text{ in some } \mathbb{R}^N \text{ for some function } \phi$$

Mercer's Theorem states the properties that K needs to satisfy to be a dot product in \mathbb{R}^N

Using one such Kernel :

$$f(x) = \beta_0 + \sum_{i=1}^n \gamma_i y_i K(x_i, x)$$

SUPPORT VECTOR

MACHINE
(Vapnik, 1990)

And feed $K(x_i, x)$ to the optimisation problem.

MOST COMMON KERNELS

LINEAR KERNEL

$$K(\underline{x}, \underline{z}) = \sum_{j=1}^p x_j z_j \quad (\text{standard inner product, SVC, linear decision surface})$$

POLYNOMIAL KERNEL

$$K(\underline{x}, \underline{z}) = (1 + \sum_{j=1}^p x_j z_j)^m \quad \text{degree } m$$

RADIAL (RBF) KERNEL

$$K(\underline{x}, \underline{z}) = \exp\left(-\gamma \sum_{j=1}^p (x_j - z_j)^2\right) \quad \gamma > 0$$

\underline{z} far from \underline{x}
 $\rightarrow K(\underline{x}, \underline{z})$ small

GAUSSIAN KERNEL

Local

Tuning parameters typically selected by cross validation.

EXAMPLE : QUADRATIC KERNEL ($m=2$)

$$\begin{aligned} K(\underline{x}, \underline{z}) &= \left(1 + \sum_{j=1}^p x_j z_j\right)^2 = \sum_{j,k=1}^p x_j x_k z_j z_k + 2 \sum_{j=1}^p x_j z_j + 1 \\ &= \sum_{j,k=1}^p (x_j z_j)(x_k z_k) + \sum_{j=1}^p (\sqrt{2} x_j)(\sqrt{2} z_j) + 1 \end{aligned}$$

So

$$\phi(\underline{x}) = \underbrace{[x_1^2, x_1 x_2, \dots, x_p^2]}_{P^2}, \underbrace{[\sqrt{2} x_1, \sqrt{2} x_2, \dots, \sqrt{2} x_p]}_P, 1$$

(x_1, \dots, x_p)

$$K(\underline{x}, \underline{z}) = \phi(\underline{x}) \cdot \phi(\underline{z}) \quad \text{with } \phi: \mathbb{R}^p \rightarrow \mathbb{R}^{p^2+p+1}$$

This kernel leads to a quadratic decision surface in the input space.

MORE THAN 2 CLASSES

Two possible approaches:

① ONE VERSUS THE REST

② ONE VERSUS ONE

① K classes \rightarrow Train K binary classifiers where the labels are "class c " versus "not class c " for $c=1, \dots, K$

\underline{x} test observation $\begin{cases} f_1(\underline{x}) \\ \vdots \\ f_K(\underline{x}) \end{cases} \rightarrow$ classify \underline{x} to the class with largest $|f_c(\underline{x})|$

② Fit $\binom{K}{2} = \frac{K(K-1)}{2}$ classifiers, one for each pair of classes

\Rightarrow classify \underline{x} to class with the majority vote among the $\frac{K(K-1)}{2}$ predictions

SVM vs LOGISTIC REGRESSION

SVC

$$\min_{\beta_0, \beta, \xi} \frac{1}{2} \|\beta\|^2$$

s.t. ① $y_i(\beta_0 + \beta^T x_i) \geq 1 - \xi_i$

② $\xi_i \geq 0$

③ $\sum_{i=1}^n \xi_i \leq C$

$\left. \begin{array}{l} \xi_i \geq 1 - y_i(\beta_0 + \beta^T x_i) \\ \xi_i \geq 0 \end{array} \right\} \rightarrow$

$$①② \Rightarrow \xi_i \geq \max\{0, 1 - y_i(\beta_0 + \beta^T x_i)\} = [1 - y_i(\beta_0 + \beta^T x_i)]_+$$

$[a]_+ = \max\{0, a\}$

Using ③:

$$\sum_{i=1}^n [1 - y_i(\beta_0 + \beta^T x_i)]_+ \leq \sum_{i=1}^n \xi_i \leq C \Rightarrow \sum_{i=1}^n [1 - y_i(\beta_0 + \beta^T x_i)]_+ \leq C$$

So the original problem can be reformulated as:

$$\min_{\beta_0, \beta} \frac{1}{2} \|\beta\|^2$$

s.t. $\sum_{i=1}^n [1 - y_i(\beta_0 + \beta^T x_i)]_+ \leq C$

Using Lagrange multipliers this becomes

$$\min_{\beta_0, \beta} \frac{1}{2} \|\beta\|^2 + \lambda \left[\sum_{i=1}^n [1 - y_i(\beta_0 + \beta^T x_i)]_+ \right]$$

$$= \min_{\beta_0, \beta} \frac{1}{2} \|\beta\|^2 + \lambda \sum_{i=1}^n [1 - y_i(\beta_0 + \beta^T x_i)]_+$$

$$= \min_{\beta_0, \beta} \sum_{i=1}^n [1 - y_i(\beta_0 + \beta^T x_i)]_+ + \frac{\lambda}{2} \|\beta\|^2$$

Loss

Ridge regression

$$\sum_{i=1}^n (y_i - \beta^T x_i)^2 + \lambda \|\beta\|^2$$

HINGE
LOSS

$$L(y, f(x)) = \begin{cases} 0 & 1 - y f(x) < 0 \iff y f(x) > 1 \\ 1 - y f(x) & 1 - y f(x) \geq 0 \iff y f(x) \leq 1 \end{cases}$$

No loss
for points
on correct
side of
margin

If on wrong side of
margin, loss is proportional to
distance from margin

Let us now consider logistic regression

OBJECTIVE FUNCTION $-\sum_{i=1}^n \log(p(c_i | x_i))$

LOSS $L(y_i, \underbrace{f(x_i)}_{\beta_0 + \beta^t x_i}) = \begin{cases} -\log p(1|x_i) & y_i = 1 \\ -\log(1-p(1|x_i)) & y_i = -1 \end{cases}$

$$y_i = 1 \quad -\log p(1|x_i) = \log\left(\frac{1}{p(1|x_i)}\right) =$$

$$= \log\left(\frac{1+e^{\beta^t x_i}}{e^{\beta^t x_i}}\right) = \log\left(e^{\cancel{-\beta^t x_i}} + 1\right)$$

$$p(1|x_i) = \frac{e^{\beta^t x_i}}{1+e^{\beta^t x_i}}$$

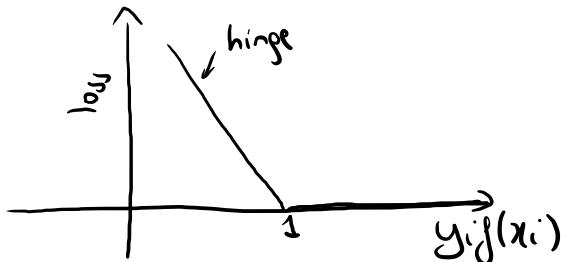
$$y_i = -1 \quad -\log(1-p(1|x_i)) = \log\left(\frac{1}{1-\frac{e^{\beta^t x_i}}{1+e^{\beta^t x_i}}}\right) = \log\left(\frac{1+e^{\cancel{\beta^t x_i}}}{e^{\beta^t x_i}}\right)$$

$$L(y_i, f(x_i)) = \log\left(1+e^{-y_i f(x_i)}\right)$$

So logistic regression in an SVM framework would be:

$$\min_{\beta_0, \beta} \sum_{i=1}^n \log(1 + e^{-y_i f(x_i)}) + \frac{\gamma}{2} \|\beta\|^2$$

LOGISTIC REGRESSION
WITH RIDGE PENALTY



In general, logistic regression and SVC perform similarly when the decision surface is linear.

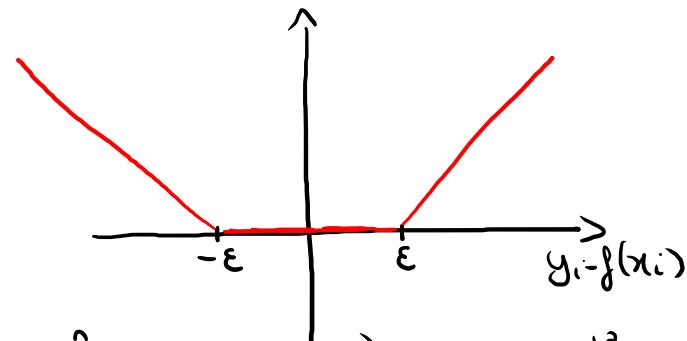
Support vector machines perform better on non-linear problems (unless you make logistic regression non-linear by adding features)

SVM in regression

The hinge loss is replaced by a ϵ -INSENSITIVE LOSS

ϵ -INSENSITIVE
LOSS

$$L_\epsilon(y_i, f(x_i)) = \begin{cases} 0 & \text{if } |y_i - f(x_i)| \leq \epsilon \\ |y_i - f(x_i)| - \epsilon & \text{if } |y_i - f(x_i)| > \epsilon \end{cases}$$



SVM
REGRESSION

$$\min_{\beta_0, \beta} \sum_{i=1}^n L_\epsilon(y_i, f(x_i)) + \frac{\lambda}{2} \|\beta\|^2$$