Let us construct $L^*(\lambda)$ for SVM:

$$L^*(\alpha_i, \lambda_i) = \min_{\omega, \omega_0} \frac{1}{2}\|\omega\|^2 + C\sum_i \xi_i$$

$$- \sum_i \alpha_i \left(y_i(\omega^T\phi(x_i) + \omega_0) + \xi_i - 1\right)$$

(assuming $\alpha_i \geq 0$, $\lambda_i \geq 0$, the RHS argument of min is strictly convex)

$$- \sum_i \lambda_i \xi_i$$

$$= \frac{1}{2}\left(\sum_i \alpha_i y_i \phi(x_i)\right)^T\left(\sum_j \alpha_j y_j \phi(x_j)\right)$$

(Substituting for $\omega$ from KKT necessary conditions)

$$+ \sum_i \xi_i \left(C - \alpha_i - \lambda_i\right)^0$$

$$- \sum_i \alpha_i \left(y_i \sum_j y_j \alpha_j \phi^T(x_j)\phi(x_i)\right)$$

$$+ \sum_i \alpha_i - \sum_i \alpha_i y_i \omega_0$$

$$= -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \phi^T(x_i) \phi(x_j)$$

$$+ \sum_i \alpha_i - w_0 \sum \cancel{\alpha_i y_i} \quad \text{○}$$

$$\boxed{L^*(\alpha_i, \lambda_i) = -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \phi^T(x_i) \phi(x_j) + \sum_i \alpha_i}$$

↓

The dual function

Our dual problem is
(Refer to class notes
of 22/01/2013)

$$\max \quad -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \phi^T(x_i) \phi(x_j)$$

$$+ \sum_i \alpha_i$$

new set of constraints ←

$$\alpha_i \geq 0$$

$$\alpha_i \geq 0 \quad \& \quad \alpha_i \leq C \quad \& \quad \sum \alpha_i y_i = 0$$

The constraints $\left( \lambda_i + \alpha_i = C \text{ 4} \right.$
$\left. \sum \alpha_i y_i = 0 \right)$
obtained while solving to obtain
the dual function $\wedge$ will also reflect in
the constraints in the dual
optimisation problem

Our final dual optimisation problem:

$$\max_{\alpha_i} \quad -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \, \phi^T(x_i) \, \phi(x_j)$$
$$+ \sum_i \alpha_i$$

$$\text{s.t} \quad C \geq \alpha_i \geq 0 \qquad [Dsum]$$

$$\sum_i \alpha_i y_i = 0$$

We hope that the solution $\alpha_i$

to the dual [DSVM] problem

recovers solution $(w, w_0)$ to
the primal SVM.

And in fact by Theorem 82
from class notes of 22/1/2013,
the duality gap $= 0$, i.e

DSVM gives us solution to
the original svm problem.

Q: Why is the dual interesting

i) Mainly for this course:
The $\phi^T(x_i)\phi(x_j)$

Let us try to write even the svm decision function in terms of the dot product.

$$f(x) = w^T \phi(x) + w_0$$

$$= \sum_i \alpha_i y_i \phi^T(x_i) \phi(x) + w_0$$

Can I write $w_0$ also in terms of the dot product?

Using the result from last class (22/1/2013):

for any $\alpha_i \in (0, c)$

$$y_i(w^T \phi(x_i) + w_0) = 1$$

$$\Rightarrow w_0 = y_i - \sum_j \alpha_j y_j \phi^T(x_i) \phi(x_j)$$

The bottom line is that if
I could compute $\phi^T(x_i)\phi(x_j)$
more directly than having
to first compute $\phi(x)$
vector & then compute

dot products, I could capture
very complex objects $(x_i)$

Let $\phi^T(x_i)\phi(x_j) = K(x_i, x_j)$

The kernel
function.

$K(x_i, x_j)$ is a kernel fn iff
$\exists \phi: X \rightarrow \mathbb{R}^m$ & $K(x_i, x_j) = \phi^T(x_i)\phi(x_j)$

In terms of the kernel function the dual svm (DSVM) problem is:

$$\max_{[\ldots \alpha_j \ldots]} \quad -\frac{1}{2} \sum_i \sum_j y_i y_j \alpha_i \alpha_j K(x_i, x_j)$$

$$+ \sum_i \alpha_i \qquad (\star)$$

$$\text{s.t} \quad \sum_i \alpha_i y_i = 0$$

$$\alpha_i \in [0, c]$$

It turns that there are efficient solvers for DSVM

(SMO algo: Sequential, Minimal Optimization)

Some guiding heuristics for picking pairs of $(\alpha_i, \alpha_j)$ for solving $(\star)$

(0) Start with all $\alpha_i = 0$ or random

(1) Pick in pairs $(\alpha_i, \alpha_j)$
since $\sum \alpha_i y_i = 0$. Thus,
you cannot pick only one
$\alpha_i$ at an iteration

(2) You can choose the $\alpha_i / \alpha_j$
in (1) based on violation of
constraints in the KKT
(see last class notes)

Eg: $\forall \; \alpha_i \in (0, c)$

$y_i(w^T \phi(x_i) + w_0) = 1$

or $\forall \alpha_i = 0$

$$y_i \left( \omega^T \phi(x_i) + w_0 \right) \geq 1$$

OR  $\neq \alpha_i = C$
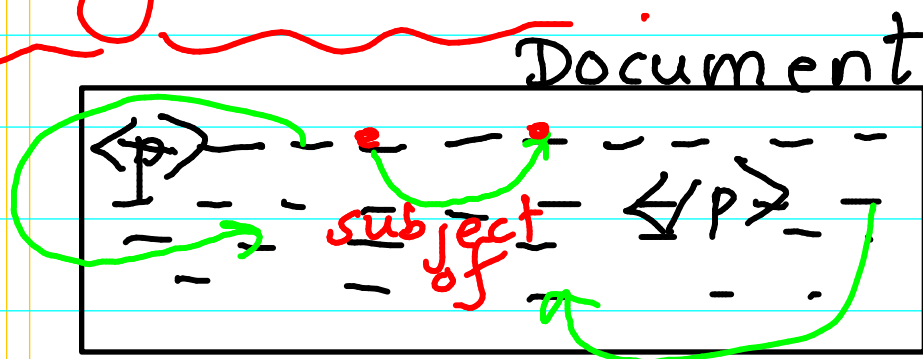
$$y_i \left( \omega^T \phi(x_i) + w_0 \right) \leq 1$$

Examples of $\phi(x)$ for which $k(x_i, x_j)$ can be computed efficiently without enumerating

$$\left[ \phi_1(x) \; \phi_2(x) \; \ldots \; \phi_m(x) \right]$$

Eg: String kernels

Subsequence kernels

Tree kernels

Graph kernels

Convolution kernels
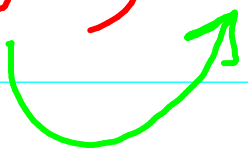
Rational kernels

First order relational kernels

Examples of kernels that enable relational learning

String kernels:

Document



Can be seen as a sequence of words

→ Can be seen as a DOM tree

↓ seen as a graph

Ram ate the apple

(subject) (verb)

In string kernel, we treat the object (document) as a sequence of tokens/words/characters

Doc D1

$w_1$ $w_2$ $w_3$

$\cdots \cdot w_t$

$(w_1 \cdots w_t$ are words$)$

Doc D2

$v_1$ $v_2$ $\cdots$

$\cdots \cdot v_r$

$(v_1 \cdots v_r$ are words$)$

A simple model:

$$\phi(D) = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

→ $i^{th}$ elements indicates if word $u_i$ occurs in $D$

A more complex model:

$$\phi(D) = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

→ element indicates if word sequence

$[u_{i_1} \; u_{i_2} \; .. \; u_{i_p}]$

occurs in $D$