1.

Support Vector Machine(SVM)

Linear Binary SVM:
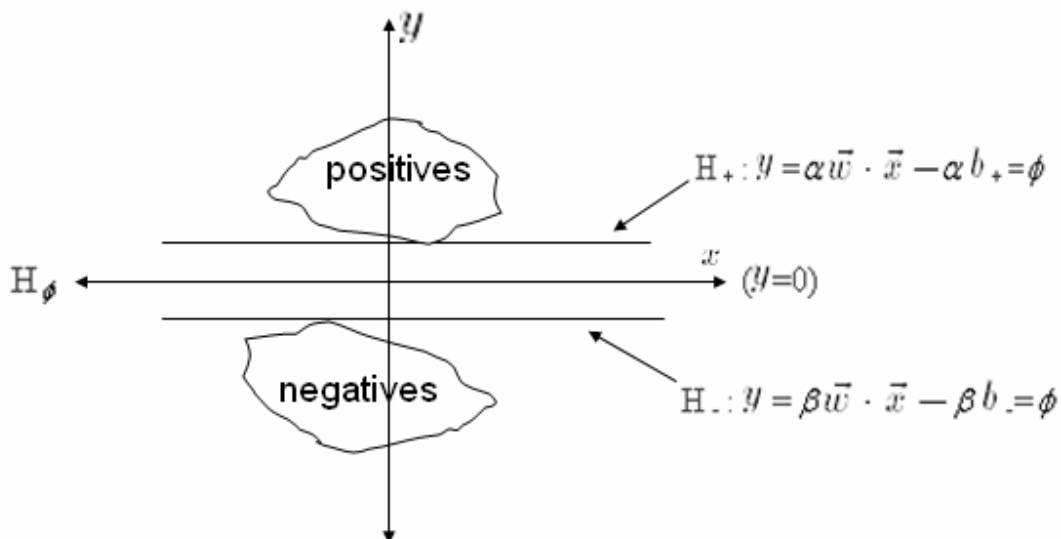
N training data "points" (feature vectors with binary labels):

$$\{\,(\bar{x}_1, y_1),(\bar{x}_2, y_2),...,(\bar{x}_n, y_n)\,\},\ \bar{x}_i \in \Re^m,\ y_i \in \{\pm 1\}$$

Assumption: The positive and negative labled data( $y_i = \pm 1$ )is sufficiently separable and "dumped"(as positives and negatives) that notions such as positive and negative data clusters, and a hyperplane separating them, are meaningful. Often this is achieved manually, through the choice of feature vector components used to represent the data instances, separability of positives and negatives is often achieved without specific choice of feature vector, however, instead leaving this to (automated) tuning at the level of the SVM kernel(to be discussed later).

Want a separating hyperplane between positives and negatives, $y_i = \pm 1$, assume full separability possible with choice of feature vector ( $f.v.$ )components:

$$H_\phi :\ y = \bar{\omega}\cdot\bar{x} - b = \phi$$



All hyperplanes parallel,
$H_\phi \parallel H_+ \parallel H\_ \rightarrow$ all have hyperplanes proportioned to $\bar{\omega}$.

2.

So, have without loss of generality (w,l,o,g):

$H_+ : \quad y = \vec{\omega} \cdot \vec{x} - b_+ = \phi$

$H_- : \quad y = \vec{\omega} \cdot \vec{x} - b_- = \phi$

Again, rescaling on $\vec{\omega}$ possible (w,l,o,g) to bring m to form:

$H_+ : \quad y = \vec{\omega} \cdot \vec{x} - b_+ = +1$

$H_- : \quad y = \vec{\omega} \cdot \vec{x} - b_- = -1$

So, for fully separable binary data $(\vec{x}_i, y_i)$:

$\vec{\omega} \cdot \vec{x}_i - b \geq +1$ for $y_i = +1$

$\vec{\omega} \cdot \vec{x}_i - b \leq -1$ for $y_i = -1$

Combined using sign trick:

$$y_i (\vec{\omega} \cdot \vec{x}_i - b) - 1 \geq 0 \quad \forall_i$$

This is the key constraint that must be satisfied for separable(binary) data.

With rescaling the separation between $H_+$ and $H_-$ becomes:
First, the boundary hyperplanes are:

$$\vec{x}_i^{(+)} \cdot \omega - b = 1; \quad \vec{x}_i^{(-)} \cdot \omega - b = -1$$

Where $\vec{x}_i^{(+)}$ are $f.v.^{'s}$ from $\vec{x}_i$ that have $y_i = +1$, and that reside on("support") the $H_+$ boundary (i.e., a support vector , "S.V."). Likewise for the $\vec{x}_i^{(-)}$.
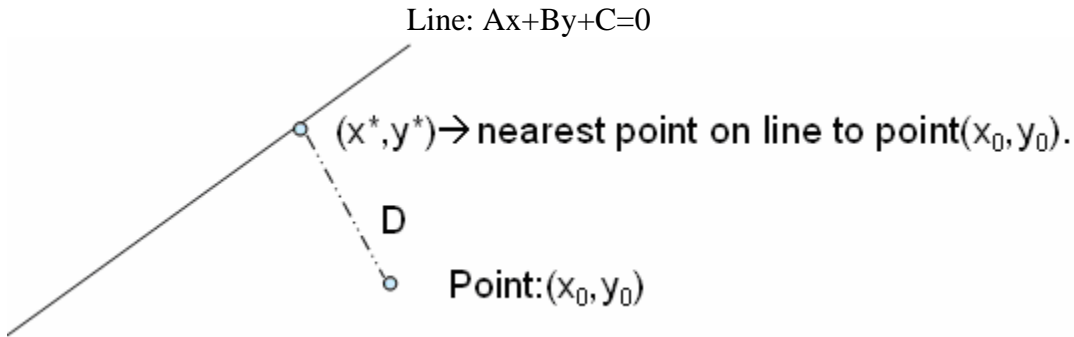
$$(\vec{x}_i^{(+)} - \vec{x}_i^{(-)}) \cdot \omega = 2$$

When $(x_i^{(+)} - x_i^{(-)})$ is perpendicular to the hyperplane the distance between the hyperplanes is given by d=$||(\vec{x}_i^{(+)} - \vec{x}_i^{(-)})|| = \dfrac{2}{||\omega||}$.     d= $2||\omega||^{-1}$

3.

A variational derivation of the result d= $\dfrac{2}{\|\omega\|}$, the distance between hyperplanes $H_+$ and $H_-$, is now shown. The variational derivation is meant to provide a refresher on methods to be used in what follows,

2-D space:

Line: Ax+By+C=0



$(x^*,y^*) \rightarrow$ nearest point on line to point$(x_0, y_0)$.

D

Point:$(x_0, y_0)$

$$D=\sqrt{(x^*-x_0)^2+(y^*-y_0)^2}$$

Want to minimize D subject to constraint $Ax^*+By^*+C = \phi$ (i.e., that the nearest point to $(x_0,y_0)$ reside on the line. This suggests the following Lagrangian formulation:

$L(x^*,y^*,\alpha)=D(x^*,y^*)+\alpha[\,Ax^*+By^*+C]$

The Lagrangian solution is obtained by minimizing L on choice of $\{x^*, y^*\}$, i.e., minimize $D(x^*, y^*)$, but subject to the constraint $Ax^*+By^*+C=0$(encapsulated in the term with the Lagrange Multiplier):

$$\frac{\partial L}{\partial \alpha} = [Ax^*+By^*+C],\ \text{requiring}\ \frac{\partial L}{\partial \alpha}=0\ \text{then restores on constraint,}$$

$$0=\frac{\partial L}{\partial x^*}=\frac{(x^*-x_0)}{D}+\alpha\,A;\ 0=\frac{\partial L}{\partial y^*}=\frac{(y^*-y_0)}{D}+\alpha\,B$$

$$\frac{(x^*-x_0)^2}{D^2}=\alpha^2 A^2\,;\ \frac{(y^*-y_0)^2}{D^2}=\alpha^2 B^2 \Rightarrow 1=\alpha^2(A^2+B^2)$$

4.

$Ax^* + By^* + C = 0 \rightarrow [Ax_0 + By_0 + C] - \alpha(A^2 + B^2)D = \phi$

$\downarrow$  $\quad\quad\quad\quad\quad\quad\quad\quad$ $\downarrow$ $\quad$ $\downarrow$

$\downarrow$ $\quad\quad\quad\quad\quad\quad\quad\quad$ Sign choice, abs maintains consisting

$\downarrow$ $\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ $\downarrow$

$Ax^* = -\alpha A^2 D + Ax_0$ $\quad\quad\quad\quad$ $D = \dfrac{Ax_0 + By_0 + C}{\sqrt{A^2 + B^2}}$

$By^* = -\alpha B^2 D + By_0$

Generalization from 2-D space to m-D space, consider distance from a point on $H_+$ to $H_\phi$:

$|Ax_0 + By_0 + C| \rightarrow |\bar{\omega} \cdot \bar{x} - b| = 1$

While the "$\bar{\omega}$" parameters in $H_\phi$ (same as in $H_+$) correspond to:

$\sqrt{A^2 + B^2} \rightarrow \sqrt{\sum_k w_k^2} = \sqrt{\bar{\omega} \cdot \bar{\omega}} = \|\bar{\omega}\|$

So, $D = \dfrac{1}{\|\bar{\omega}\|}$ from $H_+$ to $H_\phi$, twice that for $H_+$ to $H\_$:

$d = \dfrac{2}{\|\omega\|}$

The SVM approach encapsulates a Key Structured Risk Minimization (SRM) criterion when it seeks to obtain the separable solution for which "d" is the greatest. This is the solution for which the separating hyperplane is the furthest distance possible from the positive & negative supper vectors (the nearest data points). The risk assumed in using a hyperplane to separate is, intuitively, lessened if we aient that separating hyperplane to maximize its distance to the training data (least sensitivity to $s.v.^{'s}$, etc)

For Separability we have the Constraint: $y_i(\bar{\omega} \cdot \bar{x}_i - b) - 1 \geq 0 \ \forall_i$,

For SRM we have maximize $d = \dfrac{2}{\|\omega\|}$, or **minimize** $\|\omega\|^2$ (maximizing $\|\omega\|^{-2}$, instead of $\|\omega\|^{-1}$), is chosen due to simplifications in the formalism that follows---intuitively, if we max $\dfrac{1}{\|\omega\|}$ by min on $\|\omega\|$, it could just as well be done with min on $\|\omega\|^2$.

5.

Lagrangian Formulation:

$$L(\bar{\omega},b,\bar{\alpha})=\frac{1}{2}\|\omega\|^2 - \sum_i \alpha_i[y_i(\omega\cdot x_i - b)-1], \alpha_i \geq 0$$

As with the practice Lagrangian described earlier, we seek to minimize L on $\{\bar{\omega},b\}$ and to extremize(maximize in this case) L on$\{\bar{\alpha}\}$,i.e., what results due to the simultaneous minimization/ maximization is what is called the saddle-point optimization for the solution.

Note how the inequality constraints above differ from the practice (distance-to-line) problem. In the latter case, the constraint was an exact equality (Ax*+By*+C=$\phi$), and the term entering the Lagrangian was:
"$\alpha$[Ax*+By*+C]"

For which the recovery of the constraint from $\frac{\partial L}{\partial \alpha}=0$ was clear. Now, the Lagrange

Multipliers,$\alpha$, are no longer free to be positive or negative(recall the $\alpha =\pm(A^2+B^2)^{-1}$

solution before). Now the $\alpha$'s are restricted to be positive, and the term entering the Lagrangian has an overall negative in front:

"-$\alpha_i[y_i(\omega\cdot x_i - b)-1]$"

Where there are as many constraint "$i$" as there are training data points. The way to understand these constraints and their Lagrangian contributions is to directly consider the Lagrangian in an incremented saddle-point optimization:

Knows as the Karush-Kuhn-Tucker relations, i.e., the "KKT relations"
$$\downarrow$$

- If $[y_i(\omega\cdot x_i - b)-1]>\phi$ (constraint satisfied), then maximization(on $\alpha_i^{'s}$ )
  for "$\sum_i \alpha_i[y_i(\omega x_i - b)-1]$" is achieved for $\alpha_i \to \phi$ (since $\alpha_i \geq \phi$ constraint!).
- If $[y_i(\omega\cdot x_i - b)-1]=\phi$ (constraint satisfied, a support vector), then there is no constraint or $\alpha_i$.
- If $[y_i(\omega\cdot x_i - b)-1]<\phi$, then $\alpha_i \to \infty$!

6.

The last case, $[\, y_i(\omega \cdot \bar{x}_i - b) - 1\,] < \phi$, is an example of where the constraint is not satisfied. For completely separable data this case will not occur in the solution, but may occur when incrementally optimizing to achieve that solution.

As we shall see, non-separable (perfectly) data can have constraint violations in the solution. How is this managed if the Lagrangian optimization will drive the associated Lagrange multipliers to larger and larger positive values ($\alpha_i \to \infty$)?

The answer is to establish a max $\alpha$ cut off:
max($\alpha_i$)=C

Practically speaking the above is imposed for both separable and non-separable data.

When we consider the derivation of the "Dual Formalism" for L($\omega, b, \alpha$), and compare to the same Dual on the non-separable formulation--- we will find that the Duals are the same(which is very convenient). Here is the formalism w/wo separability:

$$L = \frac{1}{2}\|\omega\|^2 - \sum_i \alpha_i[y_i(\omega \cdot x_i - b) - 1], \alpha_i \geq 0,$$

[Note: $C \geq \phi \to \alpha\, C$ term in Lagrangian with $\alpha \geq \phi$]

From an implementation stand point, if nothing else, have max($\alpha_i$)=C, so, practically speaking:

$$L = \frac{1}{2}\|\omega\|^2 - \sum_i \alpha_i[y_i(\omega \cdot x_i - b) - 1], \alpha_i \geq 0, \alpha_i \leq C \quad (C - \alpha_i \geq 0)$$

The (C-$\alpha_i \geq 0$) constraint can itself be absorbed into the Lagrangian:

$$\underline{\text{minimize}} \qquad\qquad\qquad \text{interpretive (minimize) (maximize)}$$
$$\downarrow \qquad\qquad\qquad\qquad\qquad\qquad \downarrow \qquad\qquad \downarrow$$
$$L = \frac{1}{2}\|\omega\|^2 - \sum_i \alpha_i[\underline{y_i(\omega x_i - b) - 1}] + \sum \sigma_i(C - \underline{\alpha_i}), \quad \sigma_i \geq \phi, \quad \alpha_i \geq \phi$$

positive variants                    negative variants

7.

If we rewrite the Lagrangian:

$$L=\frac{1}{2}\|\omega\|^2-\sum_i\alpha_i[y_i(\omega\cdot x_i-b)-1+\sigma_i]+C\sum\sigma_i\ ,\ \sigma_i\geq 0,\alpha_i\geq 0$$

--------------------------

Corresponds to constraint:
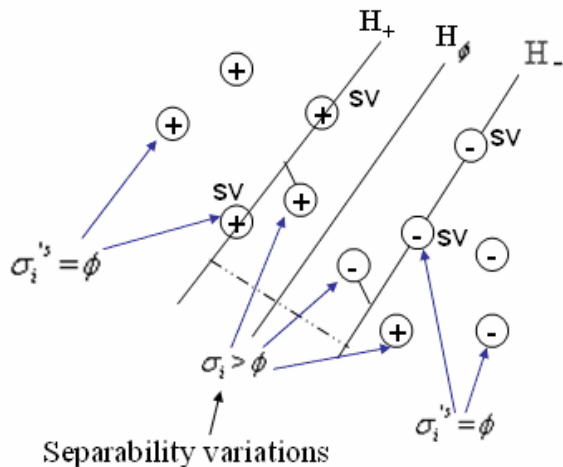$$y_i(\omega\cdot x_i-b)-1+\sigma_i\geq\phi$$
$$y_i(\omega\cdot x_i-b)\geq 1-\sigma_i\qquad\sigma_i\geq\phi$$

Back tracking: $\omega\cdot x_i-b\geq 1-\sigma_i$ for $y_i=+1$

$$\omega\cdot x_i-b\leq -1+\sigma_i\ \text{for}\ y_i=-1$$

So, the Lagrange Multiplier $\sigma_i$, introduced to deal with the max($\alpha_i$)=C constraint, can be interpreted as a "slack variable":



Separability variations

- If $C>\alpha_i,\sigma_i\rightarrow\phi$
- If $C=\alpha_i,\sigma_i$ free($\geq\phi$)
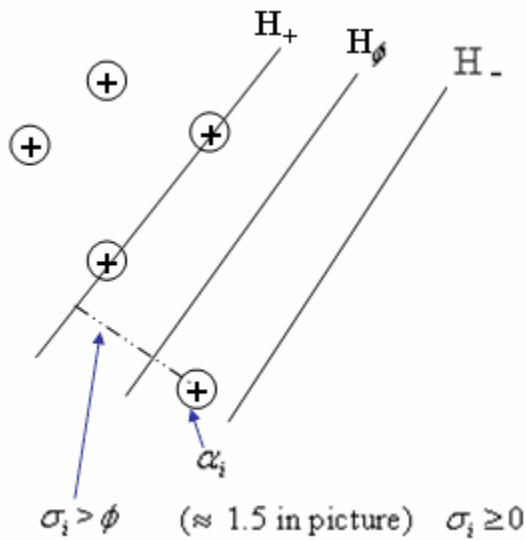- If $C<\alpha_i,\sigma_i\rightarrow\infty$

$$\downarrow$$

See Note

Note:

With the $\alpha^{'s}$ we have more control than with $[y_i(\omega\cdot x_i-b)-1]$, where previously $\alpha_i\rightarrow\infty$ resulted. Now, can avoid $C<\alpha_i$ condition by establishing initial conditions without $C<\alpha_i$ and maintaining these conditions as the Lagrangian optimization stars forward(such freedom of initialization not possible when "given" the training data{ $\bar{x}_i,y_i$ } as the starting point).

8.

To recap, a slack-variable formalism, to deal with non-perfect separability scenarios, naturally arises:



$\sigma_i > \phi$      ($\approx 1.5$ in picture)    $\sigma_i \geq 0$

If we penalize violations with a term C, then a Lagrangian modification results to

$\downarrow$            **(minimize)**

$\downarrow$                 $\sigma_i$        $(C \sum_i \sigma_i)$

The new slack constraint is:    $\sum_i \alpha_i [y_i(\omega \cdot x_i - b) - 1 + \sigma_i]$

Together:

$$L = \frac{1}{2} \| \omega \|^2 - \sum \alpha_i [y_i(\omega \cdot x_i - b) - 1 + \sigma_i] + C \sum \sigma_i \, , \, \sigma_i \geq 0, \alpha_i \geq 0$$

Exactly what we had from the max($\alpha_i$)=C constraint!

9.

Dual Calculations:

$$L = \frac{1}{2} \| \omega \|^2 - \sum_i \alpha_i [y_i(\omega \cdot x_i - b) - 1], \ \alpha_i \geq 0$$

$$L_\sigma = \frac{1}{2} \underline{\| \omega \|^2} - \sum_i \alpha_i [y_i(\omega \cdot x_i - b) - 1 + \sigma_i] + C \sum_i \sigma_i, \ \sigma_i \geq 0, \alpha_i \geq 0$$

$$\downarrow$$

$$\sum_j \omega_j^2$$

$$0 = \frac{\partial L}{\partial \omega_j} = \omega_j - \sum_i \alpha_i y_i (x_i)_j \quad \forall_j \Rightarrow \quad \underline{\bar{\omega} = \sum_i \alpha_i y_i \bar{x}_i}$$

$$0 = \frac{\partial L}{\partial b} = \sum_i \alpha_i y_i \Rightarrow \sum_i \alpha_i y_i = \phi$$

Notice that $\dfrac{\partial L_\sigma}{\partial \omega_j} = \dfrac{\partial L}{\partial \omega_j}$ and $\dfrac{\partial L_\sigma}{\partial b} = \dfrac{\partial L}{\partial b}$, so in the duality transform, where we shift to dual variables without direct reference to $\{\omega, b\}$, it will much the same:

$$\tilde{L} = \frac{1}{2} \sum_j \left( \sum_i \alpha_i y_i (x_i)_j \right)^2 - \sum_{i'} \alpha_{i'} y_{i'} \left( \sum_j \left( \sum_i \alpha_i y_i (x_i)_j \right)(x_{i'})_j \right) + \sum_i \alpha_i$$

$$= \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \sum_k (x_i)_k (x_j)_k$$

$$\boxed{\tilde{L}(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \bar{x}_i \cdot \bar{x}_j, \ \alpha_i \geq 0}$$

Where we want to find the $\alpha'^s$ that maximize $\tilde{L}(\alpha)$. Notice how in the Dual Formalism the dependence on the training data is made very clear in the $\sum_{i,j} \alpha_i \alpha_j y_i y_j \bar{x}_i \cdot \bar{x}_j$ term.

10.

For the $L_\sigma$ Dual:

$$L_\sigma \to \tilde{L}(\sigma,\alpha) = \sum_i \alpha_i - \frac{1}{2}\sum_{i,j}\alpha_i\alpha_j y_i y_j \bar{x}_i \cdot \bar{x}_j - \sum \sigma_i(\alpha_i - C), \ \sigma_i \geq 0, \alpha_i \geq 0$$

$$\boxed{\tilde{L}_\sigma(\alpha) = \sum_i \alpha_i - \frac{1}{2}\sum_{i,j}\alpha_i\alpha_j y_i y_j \bar{x}_i \cdot \bar{x}_j, \ C \geq \alpha_i \geq 0}$$

So, the duals are the same aside from the $C \geq \alpha_i$ ( $\max(\alpha) \leq C$ ) constraint, which is desirable(in some implementations) anyway.

Before moving on with the solution to $\tilde{L}_\alpha(\alpha)$, and thus the whole optimization/solution for the "decision hyperplane" (the trained SVM), let's examine further how the training data is entering the problem: Note the $\bar{x}_i \cdot \bar{x}_j$ term in the above(binary) classification Lagrangian. What does the $\bar{x}_i \cdot \bar{x}_j$ term represent?

$\bar{x}_i$ : feature vector "$i$"
$\bar{x}_i \cdot \bar{x}_j = \|\bar{x}_i\| \|\bar{x}_j\| \cos\theta_{ij}$

Suppose $\|\bar{x}_i\|=1 \ \forall_i \Rightarrow$ <u>unit hyperspherical data</u> (data points lie on unit hypersphere, where $\|x\|=1$)

$$\downarrow$$

$$\sqrt{\sum_j (x_i)_j^2} = 1 \Rightarrow \boxed{\sum_j (x_i)_j^2 = 1} \Rightarrow \|x_i\|$$

(Note: $\boxed{\sum_j |(x_i)_j| = 1} \to$ discrete probability distribution interpretation if $(x_i)_j \geq 0 \ \forall_{i,j}$ )

So, for data normalized to be unit hyperspherical: $\|\bar{x}_i\|=1 \|\bar{x}_j\| \Rightarrow \bar{x}_i \cdot \bar{x}_j = \cos\theta_{ij}$

11.

On unit hyperspherical data the interpretation of $\theta_{ij}$ is clear, it is the spherical arc angle between points $\vec{x}_i$ and $\vec{x}_j$ on its surface:

$$\vec{x}_i \cdot \vec{x}_j = \cos \theta_{ij}$$

For small angle(data in same cluster): $\cos \theta = 1 - \frac{1}{2}\theta^2 + ...$

$$\vec{x}_i \cdot \vec{x}_j \approx 1 - \frac{1}{2}\theta_{ij}^{\ 2}$$

In the Lagrangian the constant term does not matter:

$$\tilde{L}_\sigma(\alpha) \approx \sum_i \alpha_i - \frac{1}{2}\sum_{i,j}(\alpha_i y_i)(\alpha_j y_j) - \frac{1}{2}\sum_{i,j}\alpha_i \alpha_j y_i y_j (-\frac{1}{2}\theta_{ij}^{\ 2})$$

$$\underbrace{(\sum_{i,j}\alpha_i y_i)(\sum_j \alpha_j y_j)}$$

$$\downarrow$$

$$= 0 \text{ from } \frac{\partial L}{\partial b} \text{ constraint.}$$

So, for unit hyperspherical, proximate, feature vectors:

$$``\vec{x}_i \cdot \vec{x}_j" \approx -\frac{1}{2}\theta_{ij}^{\ 2} \approx -|\vec{x}_i - \vec{x}_j|^2$$

$$\uparrow \qquad\qquad \uparrow$$

Square arc length      Euclidean distance squared

So, for unit hyperspherical data, "$\vec{x}_i \cdot \vec{x}_j$" can be thought of as measuring a distance that has been regularized in some manner when the distance grows large(i.e. as angular coordinate limitation, or in exponentiation to be seen in what follows).

12.

So, for unit hyperspherical data:
$$\bar{x}_i \cdot \bar{x}_j \approx \begin{cases} -|\bar{x}_i - \bar{x}_j|^2 \text{ for } \bar{x}_i \text{ "near" } \bar{x}_j \\ \text{regularized } \{\bar{x}_i, \bar{x}_j\} \text{ expression for } \bar{x}_i \text{ not "near" } \bar{x}_j \end{cases}$$

We will see in a moment that we are free to manipulate our feature vectors by some mapping (with inverse) such that $\bar{x}_i \to \Phi(\bar{x}_i)$, which blurs the boundaries in the above example, for example, by

$$\Phi(\bar{x}_i) \cdot \Phi(\bar{x}_j) \approx \phi(-|\bar{x}_i - \bar{x}_j|^2) \text{ for positive, monotonically increasing } \phi, \text{ and for } \bar{x}_i \text{ "near" } \bar{x}_j.$$
$$\approx \text{ regularized } \phi(\{\bar{x}_i, \bar{x}_j\}) \text{ expression for } \bar{x}_i \text{ not "near" } \bar{x}_j \text{ (which may be}$$
accomplished in choice of $\phi$ implied for when $\bar{x}_i$ "near" $\bar{x}_j$ ).

One popular function of type $\phi(-|\bar{x}_i - \bar{x}_j|^2)$ is the Gaussian:

$$\phi_G(\bar{x}_i, \bar{x}_j) = \exp(-|\bar{x}_i - \bar{x}_j|^2 / 2\sigma^2)$$

For $\bar{x}_i \approx \bar{x}_j, \phi_G \approx 1 - \frac{1}{2\sigma^2}|\bar{x}_i - \bar{x}_j|^2$, and ignoring the unit constant as before we have the same form as with the hyperspherical case, except now with an extra tuning parameter(variance).

13.

As mentioned, we can map our feature vectors in a variety of ways, including lifting to a higher dimensioned feature space. This can be very powerful in identifying a separating hyperplane in a higher-dimensional mapping:

$$\Phi : \bar{x}_i \mapsto \Phi(\bar{x}_i) \text{ formerly } \bar{x}_i \cdot \bar{x}_j$$

$$\tilde{L}_\sigma(\alpha) \approx \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \Phi(\bar{x}_i) \cdot \Phi(\bar{x}_j)$$

The inner product $\Phi(\bar{x}_i) \cdot \Phi(\bar{x}_j)$ can be described as a special type of Kernel function:

$$K_{ij} = \Phi(\bar{x}_i) \cdot \Phi(\bar{x}_j)$$

Kernel functions expressible in this way satisfy what are known as Mercer's conditions(positive semidefinite K). Not all kernels satisfy Mercer's conditions, and are not describable as a mapping $\Phi$ on feature vectors. Although all kernels examined appear to satisfy Mercer's condition, this will not be taken as a critical limitation (until needed). So, w,l,o,g,

$$\tilde{L}_\sigma(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K_{ij}, 0 \le \alpha_i \le C$$

14.

As mentioned previously, the choice of kernel eliminates the need for refining a choice on feature vector mappings (beyond a certain point, such as requiring some consistent normalization on $f.v.^{'s}$ for example).

The Gaussian kernel is consistently one of the best performings:

$$K_G(\vec{x}_i, \vec{x}_j) = \exp(-|\vec{x}_i - \vec{x}_j|^2 \big/ 2\sigma^2)$$

The InnerProduct kernel(the naturally occurring one in the simple $f.v.$ case) is one of the worst performing:

$$K_{IP}(\vec{x}_i \cdot \vec{x}_j) = \vec{x}_i \cdot \vec{x}_j$$

How is the "regularization" in the Gaussian kernel better, and how might it be improved upon? The key is to see how the Gaussian kernel behaves with respect to it parameters: rewrite $K_G(\vec{y}, \vec{z}) = \exp(-\| \vec{y} - \vec{z} \|^2 / 2\sigma^2)$

$$\frac{\partial \ln K_G(\vec{y}, \vec{z})}{\partial y_K} = (y_K - z_K)\big/\sigma^2 \qquad \text{(Euclidean Distance)}$$

Clearly, the sign is important, as is a notion of difference. Suppose we generalize on this basis to

$$\frac{\partial \ln K_V(\vec{y}, \vec{z})}{\partial y_K} \quad \alpha - sign(y_K - z_K)\big/ 2\sigma^2$$

$$\text{Where sign y=1 if y>0}$$
$$\text{=-1 if y<0}$$

Notice: $$\frac{\partial(\sum_k |y_k - z_k|)^{1/2}}{\partial y_k} = \frac{sign(y_k - z_k)}{(\sum_k |y_k - z_k|)^{1/2}}$$

15.

So, with integrating factor "$1\Big/\sqrt{\sum_k | y_k - z_k |}$", we obtain an integrable form in one of the simplest ways possible. Note that now the sign convention is separated from the "notion of distance", here re-entering the kernel expression by way of the integration factor:

$$\frac{\partial \ln K_V(\bar{y}, \bar{z})}{\partial y_K} = (\frac{-1}{2\sigma^2})\,(\frac{\text{sign}(y_k - z_k)}{\sqrt{\sum_k | y_k - z_k |}})$$

$$\boxed{K_V(\bar{y}, \bar{z}) = \exp(-\sqrt{\sum_k | y_k - z_k |}\Big/ 2\sigma^2)}$$

This is usually the best performing kernel on the $L_1$ normed data considered in the channel current analysis
($L_1$ norm: $|\text{x}|_1 = \sum_k | x_k |$, a discrete prob. dist if $x_k > 0$ also)

Since $\sqrt[4]{L_1 - norm}$ can be a distance (triangle inequality, etc.), then the kernel,
$K_V = \exp(-d_v^2/2\sigma^2)$ satisfies Mercer's conditions.

Consider now the case where the notion of difference is not arithmetic but multiplicative,
i.e., based on $(1 - \frac{z_k}{y_k})$ rather than $(y_k - z_k)$ (for the Gaussian).
In doing so, we must restrict to $y_k \neq 0$ of course(thus, for all $f.v.$ components).

16.

As before, the sign of ($y_k - z_k$) is information preserved in ($1 - \dfrac{z_k}{y_k}$), but the latter is not

integrable. However, $\ln(y_k / z_k)$ also provides sign info—positive when $y_k > z_k$, etc., as before, and also includes a ratio(a multiplicative term). Which to go with? A combination seems best as this is integrable (but now need $z_k > 0$, to avoid negative terms in the log):

$$\frac{\partial \ln K_\sigma(\bar{y}, \bar{z})}{\partial y_K} = (\frac{-1}{2\sigma^2})[(1 - \frac{z_k}{y_k}) + \ln(\frac{y_k}{z_k})]$$

$$\boxed{K_\sigma(\bar{y}, \bar{z}) = \exp(-[D(y \| z) + D(z \| y)]/2\sigma^2)}$$

This is usually a close 2$^{nd}$ to the $K_V$ kernel, sometimes out performing. This kernel relates $f.v.^{'s}$ via relative entropy terms:

$$D(y \| z) = \sum_k y_k \ln(\frac{y_k}{z_k})$$

Symmetrization imposed

on relative entropy ($D_{yz} + D_{zy}$)

$\equiv$ Kullback-Lerbler Divergence,

$$\frac{\partial D(y \| z)}{\partial y_k} = \ln(\frac{y_k}{z_k}) + 1$$

$$\frac{\partial D(z \| y)}{\partial y_k} = -\frac{Z_k}{y_k}$$

A fundamental, symmetrized, information comparison between two probability

distributions. $(\dfrac{-1}{2\sigma^2}) \dfrac{\partial[D(y \| z) + D(z \| y)]}{\partial y_k} = (\dfrac{-1}{2\sigma^2})[(1 - \dfrac{z_k}{y_k}) + \ln(\dfrac{y_k}{z_k})]$

notice how the symmetrization is critical for coming together with a trivially matching term.

17.

The doubly novel aspect of the entropic kernel is that it would be the very first guess if one wanted to generalized from kernels based on exponentially regularized, square distances, to exponentially regularized, symmetrized, divergences (beginning with the most fundamental, symmetrized "relative entropy" also known as then Kullback-Leibler information divergence) .

Note that we begun with the supposition that sign was important, as was some well-behaved" notion of difference" (whether if be distance-based or <u>divergence-based</u>, etc.)

Remarkably, the entropic kernel $K_\sigma$ appears to satisfy Mercer's condition, when properly restricted to discrete probability distributions: $y_k > 0, \sum y_k = 1$. This is not established with precise mathematical proof, but is established through exhaustive numerical testing.

The Mercer test: $\sum_{i,j} k(\bar{x}_i, \bar{x}_j) C_i C_j \geq 0 \ \ \forall \bar{C} \in \Re^m$ (for positive semidefinite K)

18.

To Recap:

Seeking the $\bar{\alpha}$ that will maximize the following Larangian:

$$\tilde{L}_\sigma(\alpha) = \sum_i \alpha_i - \frac{1}{2}\sum_{i,j}\alpha_i\alpha_j y_i y_j K_{ij}, 0 \le \alpha_i \le C$$

<u>KKT Relations</u>: Let $f(\bar{x}_i) = \omega \cdot \bar{x}_i - b$

$x_i \cdot x_j \Rightarrow K_{ij}$

Have solution on constraints when:
- $\alpha_i = \phi \Leftrightarrow y_i f(\bar{x}_i) \ge 1$
- $0 < \alpha_i < C \Leftrightarrow y_i f(\bar{x}_i) = 1$
- $\alpha_i = C \Leftrightarrow y_i f(\bar{x}_i) \le 1$

Where:

$$f(\bar{x}_i) = \omega \cdot \bar{x}_i - b = \sum_j \alpha_j y_j (\bar{x}_j \cdot \bar{x}_i) - b$$

$$\downarrow \text{kernel generalization}( x \cdot y \rightarrow K(x,y))$$

$$f(\bar{x}_i) = \sum_j \alpha_j y_j K(\bar{x}_j \cdot \bar{x}_i) - b$$

Initialization:

Since $\sum_i \alpha_i y_i = \phi$, choose $\alpha_i^+ = \dfrac{1}{N^+}, \alpha_i^- = \dfrac{1}{N^-}$

Now to consider the solution $\bar{\alpha}$ by sequential Minimal Optimization(SMO), where successive pairs of $\alpha^{'s}$ are selected for optimization!

$\tilde{L}_\sigma(\bar{\alpha}) = \tilde{L}_\sigma(\alpha_1, \alpha_2, \alpha_3, ... \alpha_n)$ then "freeze" all but $\alpha_1$ and $\alpha_2$ in variational optimization (each pair of $\alpha^{'s}$ selected above could simply be relabeled as $\alpha_1$ and $\alpha_2$, thus use the same derivation to follow).

19.

SMO

$$\tilde{L}_\sigma(\vec{\alpha}) = \tilde{L}_\sigma(\alpha_1, \alpha_2; \alpha_3, \dots \alpha_m); \text{ Note: } \boxed{s = y_1 y_2}$$

$$= (\alpha_1 + \alpha_2 + \sum_{i \geq 3} \alpha_i) - \frac{1}{2}(\alpha_1^2 k_{11} + \alpha_2^2 k_{22} + 2\alpha_1 \alpha_2 s k_{12})$$

$$- \frac{1}{2}(2\alpha_1 y_1 \sum_{j \geq 3} \alpha_j y_j k_{1j} + 2\alpha_2 y_2 \sum_{j \geq 3} \alpha_j y_j k_{2j}) - \frac{1}{2} \sum_{j \geq 3, i \geq 3} \alpha_i \alpha_j y_i y_j k_{ij}$$

Let $v_i = \sum_{j \geq 3} \alpha_j y_j k_{ij} = \vec{\omega} \cdot \vec{x}_i - \alpha_1 y_1 k_{i1} - \alpha_2 y_2 k_{i2}$

$$\tilde{L}_\sigma(\vec{\alpha}) = \alpha_1 + \alpha_2 - \frac{1}{2}(\alpha_1^2 k_{11} + \alpha_2^2 k_{22} + 2\alpha_1 \alpha_2 s k_{12}) - \alpha_1 y_1 v_1 - \alpha_2 y_2 v_2$$

$$\underbrace{+ \sum_{i \geq 3} \alpha_i - \frac{1}{2} \sum_{j \geq 3, i \geq 3} \alpha_i \alpha_j y_i y_j k_{ij}}_{\{\alpha_1, \alpha_2\} \text{independent term s}}$$

Now consider variational parameters other than $\{\alpha_1, \alpha_2\}$ to be fixed in the $\{\alpha_1, \alpha_2\}$ variational optimization. Furthermore:

$$\sum_i y_i \alpha_i = 0 \Rightarrow y_1 \alpha_1 + y_2 \alpha_2 = -\sum_{i \geq 3} y_i \alpha_i$$

$$(\alpha_1 + s\alpha_2) = \gamma \qquad (\gamma = -y \sum_{i \geq 3} y_i \alpha_i)$$

$$\uparrow$$

does not depend on $\alpha_1$ or $\alpha_2$

$$\alpha_1 = \gamma - s\alpha_2$$

20.

$$L(\alpha_2;\alpha_3,...\alpha_m) \quad = (\gamma - s\alpha_2) + \alpha_2 - \frac{1}{2}((\gamma - s\alpha_2)^2 k_{11} + \alpha_2^2 k_{22} + 2\alpha_2(\gamma - s\alpha_2)sk_{12})$$

$$(\alpha_1 = \gamma - s\alpha_2) \quad\quad - (\gamma - s\alpha_2)y_1v_1 - \alpha_2 y_2 v_2 + [\text{terms independent of } \{\alpha_1,\alpha_2\}]$$

$$0 = \frac{\partial L}{\partial \alpha_2} = (1-s) - \alpha_2^{*}(k_{11} + k_{22} - 2k_{12}) + s\gamma k_{11} - s\gamma k_{12} + sy_1v_1 - y_2v_2$$

$$\text{Definition: } -\eta = k_{11} + k_{22} - 2k_{12}$$

new $\alpha$, the optimization solution, $\quad\quad$ old $\alpha's$

$$(-\eta)\alpha_2^* = (1-s) + s(k_{11} - k_{12})\gamma + y_2(\bar{\omega}\cdot\bar{x}_1 - \alpha_1 y_1 k_{11} - \alpha_2 y_2 k_{12})$$
$$- y_2(\bar{\omega}\cdot\bar{x}_2 - \alpha_1 y_1 k_{21} - \alpha_2 y_2 k_{22})$$

$$(-\eta)\alpha_2^* = (1-s) + s(k_{11} - k_{12})\gamma + y_2\bar{\omega}\cdot(\bar{x}_1 - \bar{x}_2) - (\gamma - s\alpha_2)sk_{11} - \alpha_2 k_{12} + (\gamma - s\alpha_2)sk_{21} + \alpha_2 k_{22}$$

$$= (1-s) + y_2[\bar{\omega}\cdot x_1 - y_1 - (\bar{\omega}\cdot x_2 - y_2) + (y_1 - y_2)] + \alpha_2 \underbrace{(k_{11} + k_{22} - 2k_{j2})}_{(-\eta)}$$

$$\boxed{\alpha_2^* = \alpha_2 - y_2[\bar{\omega}\cdot x_1 - y_1 - (\bar{\omega}\cdot x_2 - y_2)]/\eta}$$

This gives an analytic calculation for the optimal new $\alpha_2$. ($\alpha_2^*$).
The problem is that the solution to the optimum may take $\alpha_2$ out of range, i.e. $\alpha_2 > C$ not allowed. If such is indicated, then the procedure is to "clip" the $\alpha_2$ update to the boundary value ($\alpha_2 = C$):  $\quad \alpha_2^* \to \alpha_2^*|_{clipped}$

$$\alpha_1^* = \alpha_1 + s(\alpha_2 - \alpha_2^*|_{clipped})$$

21.

There are some subtleties in how to manage the clipping with constraints on both $\alpha_1$ and $\alpha_2$:

Suppose varying $\{\alpha_1, \alpha_2\}$ and $y_1 \neq y_2$, and that $\boxed{\alpha_2 {}^* = \alpha_2 + \Delta\alpha_2}$

$$\sum y_i \alpha_i = \phi \Rightarrow \Delta\alpha_1 = \Delta\alpha_2$$

Maximum $(\Delta\alpha_2) = \alpha_2 + (C - \alpha_2) = C$     |     $\Delta\alpha_2 = C - \alpha_2$ covers rest of distance

to

         or

         $= \alpha_2 + (C - \alpha_1) = C + \alpha_2 - \alpha_1$     $\alpha_2 \leq C$ boundary

For $y_1 \neq y_2$

Max $(\Delta\alpha_2) = $ Min $\{ C - \alpha_2, C - \alpha_1 \}$     $\Delta\alpha_1 = C - \alpha_1$ covers rest of distance to

Min $(\Delta\alpha_2) = $ Max $\{ -\alpha_2, -\alpha_1 \}$     $\alpha_1 \leq C$ boundary and since $\Delta\alpha_1 = \Delta\alpha_2$

         are tied together, this limits $\Delta\alpha_2 = C - \alpha_1$

For $y_1 = y_2 \Rightarrow \Delta\alpha_1 = -\Delta\alpha_2$

Max $(\Delta\alpha_2) = $ Min $\{ (C - \alpha_2), \alpha_1 \}$

Min $(\Delta\alpha_2) = $ Max $\{ -\alpha_2, \alpha_1 - C \}$

SMO allows us to reduce the computational problem significantly due to the exact analytical solutions obtained. Further, subtleties arise upon implementing in actual code, as will be discussed in the cluster, Perl code comments to follow.