Minimize $f_0(x)$ [ Convex Optimization ]

Subject to constraints (i.e. $x$ should be such that the below conditions hold)

$$f_i(x) \le 0 \qquad i = 1, \dots, m \qquad (*)$$
$$h_i(x) = 0 \qquad i = 1, \dots, p$$

Consider the Lagrangian

$$\mathcal{L}(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^{m} \lambda_i f_i(x) + \sum_{i=1}^{p} \nu_i h_i(x)$$

This gives us the Lagrange dual function:

$$g(\lambda, \nu) = \inf_{x \in D} \mathcal{L}(x, \lambda, \nu)$$

domain of $f_0(x)$ under constraints

which gives us a lower bound on the minimum of $(*)$

$$g(\lambda, \nu) \le f_0(x^*)$$

whenever $\lambda \succeq 0$ (easy to prove)

So, we now replace our original difficult minimization problem $(*)$ with an easier maximization of this lower bound to get as close as possible to the the minimum value $f_0(x^*)$. This max. problem is called the Lagrange dual problem:

$$\text{maximize } g(\lambda, \nu)$$
$$\text{subject to } \lambda \succeq 0 \qquad (**)$$

$\Rightarrow$ This is a convex optimization problem!

($x \succeq 0$ simply means all components of vector $x$ are $\ge 0$)

The optimal solution $(\lambda, v^*)$ is __dual optimal__

Any pair $(\lambda, v)$ s.t. $\lambda \geq 0$ and $g(\lambda, v) > -\infty$ is __dual feasible__

As we ~~stated~~ stated above (easily proveable), __weak duality__ always holds:

$$g(\lambda^*, v^*) \leq f_0(x^*)$$

But sometimes, __strong duality__ holds!

$$g(\lambda^*, v^*) = f_0(x^*)$$

This holds whenever __constraint qualifications__ are satisfied. One such example is:

① Primal problem is convex, i.e. $h_q(x) = A_q x - b_q = 0$, (equality constraints ~~are~~ affine)   $p = 1$

② __Slater's condition__ holds! there exists some (strictly feasible) point $\tilde{x}$ s.t. $f_i(\tilde{x}) < 0$ $\forall i$ and $A\tilde{x} = b$

If the objective $f_0$ and constraint $f_i$, $h_i$ functions are differentiable, ~~and~~ Slater's condition holds, and strong duality holds, then the __KKT conditions__ are necessary and sufficient for global optimality

- $f_i(x) \leq 0$, $h_i(x) = 0$, $\lambda_i \geq 0$ ~~/////////////~~
- $\nabla f_0(x) + \sum_{j=1}^{m} \lambda_i \nabla f_i(x) + \sum_{i=1}^{p} v_i \nabla h_i(x) = 0$

$\left( \begin{array}{l} \text{i.e. if you solve for } x \\ \text{such that the KKT conditions} \\ \text{hold, then you are at the global} \\ \text{optimum} \end{array} \right)$

- $\underline{\lambda_i f_i(x) = 0}$

$\quad\quad \hookrightarrow$ this condition is called __complementary slackness__ and it follows from strong duality:

$$f_0(x^*) \underset{\uparrow}{=} g(\lambda^*, v^*) = \inf_{x \in D} \left( f_0(x) + \sum \lambda_i^* f_i(x) + \sum v_i^* h_i(x) \right) \leq f_0(x^*) + \sum \lambda_i^* f_i(x^*) + \sum v_i^* h_i(x^*)$$

$$\therefore \sum_{i=1}^{m} \lambda_i^* f_i(x^*) = 0 \iff \begin{cases} \lambda_i^* > 0 \implies f_i(x) = 0 \\ f_i(x^*) < 0 \implies \lambda_i^* = 0 \end{cases}$$

<u>Representer Theorem</u>: Suppose we ~~have~~ have a set of data points $\{(x_i, y_i)\}_{i=1}^{N}$ and we want to find the function/ ~~input~~ input-output mapping $f(\cdot)$ that minimizes the loss function:

$$f^* = \underset{f \in \mathcal{H}}{\arg\min} \; L_y\big(f(x_1), \ldots, f(x_N)\big) + \Omega\left(\|f\|_{\mathcal{H}}^2\right)$$

where $\Omega(\cdot)$ is non-decreasing and $y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$ parameterizes $L_y(\cdot)$. Note that $L_y$ depends on $x_i$'s only via $f(x_i)$. For example, in ridge regression $L_y(f(x_1), \ldots, f(x_N)) = \frac{1}{2}\sum_{i=1}^{N}(y_i - f(x_i))^2$ and $\Omega(\|f\|_{\mathcal{H}}^2) = \lambda\|f\|_{\mathcal{H}}^2$. The theorem now tells us that a solution to this minimization takes the form:

$$f^* = \sum_{i=1}^{N} \alpha_i K(x_i, \cdot)$$

if $\Omega(\cdot)$ is strictly increasing.

<u>Pf.</u> Let $f^* = f_s + f_\perp$, where $f_s$ is the projection of $f^*$ onto the subspace spanned by $\{K(x_i, \cdot)\}_{i=1}^{N}$ and $f_\perp$ is the orthogonal error relative to $f$.

first note that

$$L_y\left(f(x_1), \ldots, f(x_N)\right)$$

$$= L_y\left(\langle f, K(x_1, \cdot)\rangle, \ldots, \langle f, K(x_N, \cdot)\rangle\right)$$

$$= L_y\left(\langle f_s + f_\perp, K(x_1, \cdot)\rangle, \ldots, \langle f_s + f_\perp, K(x_N, \cdot)\rangle\right)$$

$$= L_y\left(\langle f_s, K(x_1, \cdot)\rangle, \ldots\right)$$

$$= L_y\left(f_s(x_1), \ldots, f_s(x_N)\right)$$

So, minimizing $L_y$ w.r.t. $f_s$ is the same as minimizing w.r.t. $f$: we can forget $f_\perp$ without losing anything.

Now note that $f_s$ is in fact the minimum of $\Omega\left(\|f\|_{\mathcal{H}}^2\right)$ if it is ~~monotonically~~ non-decreasing, since in this case,

$$\Omega\left(\|f\|_{\mathcal{H}}^2\right) = \Omega\left(\|f_s\|_{\mathcal{H}}^2 + \|f_\perp\|_{\mathcal{H}}^2\right) \geq \Omega\left(\|f_s\|_{\mathcal{H}}^2\right)$$

Thus, this component is minimized when $\|f_\perp\|_{\mathcal{H}} = 0$, leaving the unique (only unique whenever $\Omega(\cdot)$ is strictly increasing) solution

$$f = f_s = \sum_{i=1}^{N} \alpha_i K(x_i, \cdot)$$

# Support Vector Classification

The problem is to find a hyperplane that separates the data correctly according to some classification criteria. Formally, what we want is a hyperplane such that the scalar projection $w^Tx_i$ of all data points onto the direction $w$ perpendicular to it gives us the correct classification:



$$y = \text{sign}\left(w^Tx_i + b\right)$$

we can find the best such hyperplane by maximizing the minimum distance b/w it and each class $(y=+1, y=-1)$, i.e. maximizing the **margin**. We can compute this by considering a pair of points of different classes $x^+, x^-$ lying on each margin: these will be at the minimum distance from the hyperplane, which

~~We impose that this minimum distance be 1, measured by the scalar projection onto w:~~
~~$w^Tx^+ + b = \min(w^Tx_i + b) = 1 \quad \forall x_i : y_i = +1$~~
~~$w^Tx^- + b = \max(w^Tx_i + b) = -1 \quad \forall$~~

That can be computed via $\dfrac{x^{+T}W}{\|w\|}, \dfrac{x^{-T}v}{\|w\|}$

Since $x^+$ is of class $y=+1$ and $x^-$ of class $y=-1$, we know that
$w^Tx^+ + b \geq 0, \quad w^Tx^- + b < 0.$

In fact, ~~we don't want~~ to ensure / we are going to enforce that accuracy of our classifier ~~~~ choice of $v, b$ such

$$w^T x_i + b \geq 1 \quad \forall i: y_i = 1 \quad \text{and} \quad w^T x_i + b \leq -1 \quad \forall i: y_i = -1$$

~~where the inequalities~~ become equalities for points on the margin ~~that are closest~~ to the hyperplane.

Maximizing the minimum distance b/w classes and the hyperplane thus entails maximizing the distance b/w margins which we now have is

$$\frac{x^{+T} w}{\|w\|} - \frac{x_i^T w}{\|w\|} = \frac{(1-b) - (-1-b)}{\|w\|} = \frac{2}{\|w\|}$$

~~Our~~ job has thus become solving the following optimization problem:

$$\text{maximize} \quad \frac{2}{\|w\|} \quad \text{subject to} \quad w^T x_i + b \begin{cases} \geq 1 & \forall i: y_i = +1 \\ \leq -1 & \forall i: y_i = -1 \end{cases}$$

which can be rewritten as

$$\min_{w,b} \frac{1}{2}\|w\|^2 \quad \text{subject to} \quad y_i(w^T x_i + b) \geq 1$$

However, it will rarely be possible to find a hyperplane that perfectly separates the ~~two~~ classes, so we soften the ~~objective~~ constraint and modify our objective to include a trade-off (controlled by C) with errors (i.e. data points within the margins or on the wrong side of the hyperplane):

$$\min_{w,b} \left( \frac{1}{2}\|w\|^2 + C \sum_{i=1}^{N} \xi_i \right) \quad \text{subject to} \quad y_i(w^T x_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0$$

This gives us the following Lagrangian:

$$\mathcal{L}(w, b, \alpha, \lambda, \xi) = \frac{1}{2}\|w\|^2 + C \sum_{i=1}^{N} \xi_i + \sum_{i=1}^{N} \alpha_i \left(1 - (w^T x_i + b) y_i - \xi_i \right) + \sum_{i=1}^{N} \lambda_i (-\xi_i)$$

Noting that each of our constraints $f_i(x) = 1 - \xi_i - (w^T x_i + b) y_i \leq 0$, $g_i(\xi_i) = -\xi_i \leq 0$ are convex, and that there always exists some $x, \xi$ that satisfies them (i.e. Slater's condition holds), we have that strong duality holds. Therefore, we need only solve for the KKT conditions to get the global optimum.

1) $z_i \geq 0$, $\alpha_i \leq 0$

2) $\frac{\partial \mathcal{L}}{\partial w} = w - \sum_{i=1}^{N} \alpha_i y_i x_i = 0 \implies \underline{w = \sum_{i=1}^{N} \alpha_i y_i x_i}$

$\frac{\partial \mathcal{L}}{\partial b} = \sum_{i=1}^{N} \alpha_i y_i = 0$

$\frac{\partial \mathcal{L}}{\partial \zeta_i} = c - \alpha_i - z_i = 0 \iff \underline{\alpha_i = c - z_i}$

$\implies \underline{\alpha_i \leq c}$   since $z_i \geq 0$

3) (complementary slackness)

For $\alpha_i = c$ $(\neq 0)$,      $x_i$ lies inside
                                    ($x_i$'s within the margins)
$z_i = 0 \implies \zeta_i \geq 0$

$1 - (w^T x_i + b) y_i - \zeta_i = 0 \iff y_i(w^T x_i + b) = 1 - \zeta_i$

For $0 < \alpha_i < c$,                $\rightarrow (x_i$ lies on margin$)$
$z_i > 0 \implies \zeta_i = 0$                           ↑
as in first case, $y_i(w^T x_i + b) = 1 - \zeta_i = 1$

For $\alpha_i = 0$
$z_i > 0 \implies \zeta_i = 0 \longrightarrow$ correctly
$y_i(w^T x_i + b) \geq 1$   $(x_i$'s $\checkmark$ outside the margins$)$

In other words, we find that our solution for $\alpha$ is such that
   - it is <u>sparse</u>: (only points on the margin or w/ ~~too~~ large error
      (i.e. inside the margins) have $\alpha_i \geq 0$

   - only <u>these points</u> contribute to the ~~support vector~~ $w = \sum_i \alpha_i y_i x_i$
   - the contribution of ~~error~~ large error $x_i$'s is bounded by $c$

thus, these are called the <u>support vectors</u>

We can now solve for the support vector by maximizing the dual $g(\alpha)$ with respect to $\alpha$. We first express the full dual $g(\alpha, \lambda)$ in terms of just $\alpha$, which we can do given our KKT conditions we derived above:

$$g(\alpha, \lambda) = \frac{1}{2}\|w\|^2 + C\sum_i \xi_i + \sum_i \alpha_i \left(1 - y_i(w^T x_i + b) - \xi_i\right) + \sum_i \lambda_i(-\xi_i)$$

$$= \frac{1}{2}\sum_{ij} \alpha_i \alpha_j y_i y_j x_i^T x_j + C\sum_i \xi_i + \sum_i \alpha_i - \sum_i \alpha_i \xi_i$$

$$= \sum_i \alpha_i y_i x_i \sum_j \alpha_j y_j x_j^T - b\underbrace{\sum_i \alpha_i y_i}_{=0} - \sum_i (C - \alpha_i)\xi_i$$

$$= -\frac{1}{2}\sum_{ij} \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum_i \alpha_i + \cancel{C\sum_i \xi_i} - \cancel{\sum_i (C - \alpha_i)\xi_i}$$

$$= g(\alpha)$$

We now simply minimize $g(\alpha)$ subject to the constraints

$$0 \leq \alpha_i \leq C$$
$$\sum_i \alpha_i y_i = 0$$

which is a quadratic program. The resulting solution then gives us the support vector $w$ by our equation derived above. We get $b$ by solving the equation $y_i(w^T x_i + b) = 1$ for an $x_i$ on the margin or by averaging the solutions for all $x_i$ on the margins.

## $\nu$-SVM

We can also give an alternative formulation of the problem that yields more interpretable parameters (as opposed to $C$, which is rather opaque). The following formulation is called $\nu$-SVM:

$$\min_{w, \rho, \xi} \left(\frac{1}{2}\|w\|^2 - \nu\rho + \frac{1}{N}\sum_{i=1}^N \xi_i\right) \quad \text{subject to} \quad \begin{array}{c} \rho \geq 0 \\ \xi_i \geq 0 \\ y_i(w^T x_i) \geq \rho - \xi_i \end{array}$$

where we have dropped the offset $b$ purely for simplicity.

the resulting Lagrangian is:

$$\mathcal{L}(w, v, \rho, \{\xi\}, \alpha, \lambda, \gamma)$$

$$= \frac{1}{2}\|w\|^2 - v\rho + \frac{1}{N}\sum_{i=1}^{N}\xi_i + \sum_{i=1}^{N}\alpha_i\left(\rho - y_i w^T x_i - \xi_i\right) + \sum_{i=1}^{N}\lambda_i(-\xi_i) + \ldots$$

here we can interpret the new parameter $\rho$ as the margin width we want to optimize, along with the support vector $w$ and the errors $\xi_i$. We now follow the same exercise as above, first writing out the KKT conditions after noting that again strong duality holds and then writing out the dual function:

1) $\alpha_i \geq 0,\ \lambda_i \geq 0,\ \gamma \geq 0$

2) $$\frac{\partial \mathcal{L}}{\partial W} = w - \sum_{i=1}^{N}\alpha_i y_i x_i = 0 \iff \underline{\underline{w = \sum_{i=1}^{N}\alpha_i x_i y_i}}$$

$$\frac{\partial \mathcal{L}}{\partial \xi_i} = \frac{1}{N} - \alpha_i - \lambda_i = 0 \iff \underline{\underline{\alpha_i + \lambda_i = \frac{1}{N}}}$$

$$\frac{\partial \mathcal{L}}{\partial \rho} = -v + \sum_{i=1}^{N}\alpha_i - \gamma = 0 \iff \underline{\underline{v = \sum_{i=1}^{N}\alpha_i - \gamma}}$$

3) Complementary slackness. Lets assume that $\rho > 0$ to consider only non-trivial cases w.r.t. our new parameter $v$. By complementary slackness, this implies that $\gamma = 0$, which implies that $v = \sum_{i=1}^{N}\alpha_i$. Now we consider two cases for $\xi_i$:

For $\xi_i > 0$:

$$\Rightarrow \lambda_i = 0 \iff \alpha_i = \frac{1}{N}$$

then, for all such points $N(\alpha)$

$$\sum_{i\in N(\alpha)}\alpha_i = \frac{|N(\alpha)|}{N} \leq \sum_{i=1}^{N}\alpha_i = v$$

Noting that $N(\alpha)$ is the set of all points that fall inside the margins, we can interpret $v$ as an upper bound on the number of such 'errors'.

For $\xi_i = 0$,

$$\lambda_i > 0 \implies \alpha_i < \frac{1}{N} \qquad \sum_{i \in N(\alpha)} \alpha_i + \sum_{i \in M(\alpha)} \alpha_i < \frac{|N(\alpha)| + |M(\alpha)|}{N} \leq \nu$$

~~RRRR10RRRRRRRR~~ $\implies \nu$ is upper bound on total # of support vectors w/ non-zero weight

Let $M(\alpha)$ be the set of points such that $0 < \alpha_i < \frac{1}{N}$, i.e. the points with $\xi_i = 0$ that still contribute to ~~the weight~~ $w$ (i.e. $\alpha_i \neq 0$)

The dual function is then:

$$g(\alpha) = \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j - \nu\rho + \frac{1}{N} \sum_i \xi_i + \sum_i \alpha_i \rho - \sum_i \alpha_i \rho_0 - \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j$$

$$- * \sum_i \left(\frac{1}{N} - \alpha_i\right) \xi_i + \left(\nu - \sum_i \alpha_i\right)\rho \quad = -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j$$

So we now maximize $g(\alpha)$ subject to: $\sum_{i=1}^{n} \alpha_i \geq \nu$

$$0 \leq \alpha_i \leq \frac{1}{N}$$

## Kernelized SVM

We can easily accomodate a kernelized solution to the problem by recognizing the form* of the objective function being minimized and invoking the representer theorem, telling us that $w = \sum_{i=1}^{N} \beta_i K(x_i, \cdot)$. We can thus interpret the minimization of $\|w\|_{\mathcal{H}}^2$ (i.e. the maximization of the margin) as enforcing smoothness of the function $w \in \mathcal{H}$.

Our objective function in terms of $\xi_i$ thus becomes (again dropping $b$ for simplicity)

$$\min_{\beta, \xi} \left(\frac{1}{2} \beta^T K \beta + C \sum_{i=1}^{n} \xi_i\right) \quad \text{subject to} \quad \xi_i \geq 0$$
$$y_i \sum_{j=1}^{n} \beta_j K(x_j, x_i) \geq 1 - \xi_i$$

Since $K$ is positive definite, this objective is convex and strong duality holds, giving the dual function

$$g(\alpha) = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(x_i, x_j),$$

which we maximize subject to $0 \leq \alpha_i \leq C$.

* In fact, to see this we need to put our objective in the form
$$\frac{1}{2} \|w\|_{\mathcal{H}}^2 + C \sum_{i=1}^{n} [1 - y_i \langle w, K(x_i, \cdot) \rangle_{\mathcal{H}}]_+ \quad \leftarrow \text{~~regularization~~} \quad \text{to invoke the representer theorem.}$$
This is equivalent to the ~~the~~ form in terms of $\xi_i$, just harder to minimize b/c of the non-linearity