# Support Vector Machines (SVMs)

Ref: Chapter 5 of FoML book by Mohri et al.

**Setting:** Two-class classification problem

Input space $X \subseteq \mathbb{R}^d$, class labels $Y \in \{+1, -1\}$

Target function $f : X \to Y$

Dataset $S = \{ (x_i, y_i), i = 2) \ldots m \}$, $y_i = f(x_i)$

$\quad \quad \quad \quad \quad x_i$

$\quad \to$ The samples drawn from some unknown distribution "D" in an iid fashion.

$$R(h) = \underset{x \sim D}{P} (h(x) \neq f(x)) \text{ is the}$$

generalization error of hypothesis $h$.

**Goal:** Minimize the generalization error

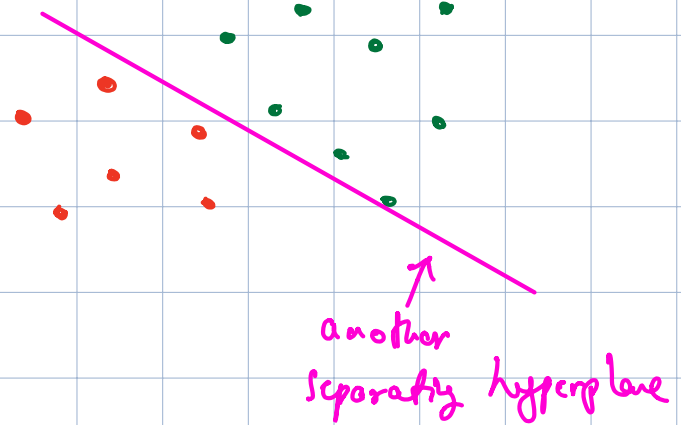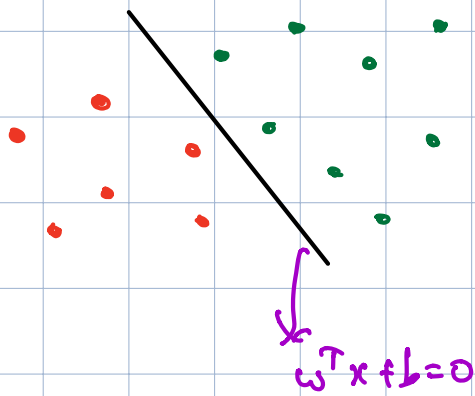The hypothesis set we consider is

$$H = \{ x \to \text{sign}(w^T x + b) : w \in \mathbb{R}^d, b \in \mathbb{R} \}$$

" Linear classification problem".

$w^T x + b = 0$ is a hyperplane & a hypothesis

$$h: x \to \text{sgn}(w^T x + b)$$



$w^T x + b = 0$

another separating hyperplane

## Linearly Separable dataset:

A dataset $S = \{(x_i, y_i), i = 1 \cdots m\}$ is linearly separable if $\exists\, w, b$ s.t.

$$w^T x_i + b > 0 \quad \forall i \text{ with } y_i = +1, \text{ and}$$
$$w^T x_i + b < 0 \quad \forall i \text{ with } y_i = -1. \qquad (\divideontimes)$$

Note: $(\divideontimes)$ is equivalent to saying $\exists\, c > 0$ s.t.

$$w^T x_i + b > c \quad \forall i,\ y_i = +1$$
$$w^T x_i + b < -c \quad \forall i,\ y_i = -1$$

> Since we have a finite # of points in $S$

$\Rightarrow$ we can scale $w, b$ to have

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

(or) equivalently, $\forall i,$ 
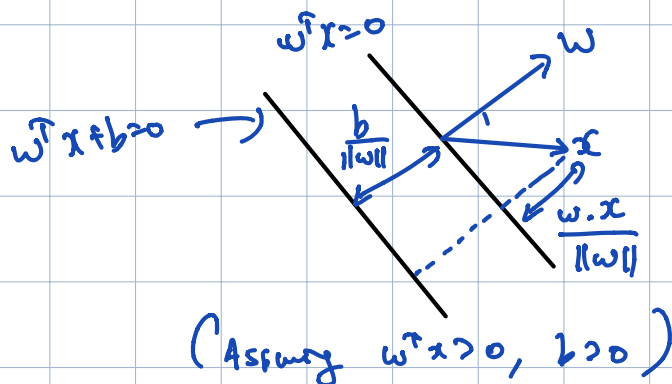$$y_i (w^T x_i + b) \geq 1 \quad \leftarrow \text{ Condition for linear separability}$$

"Assume $S$ is linearly separable".

## Concept of margin:

Margin: $\ell_h(x)$ of a linear classifier $h: x \to w^T x + b$ at the point $x$ is

$$\ell_h(x) = \frac{|w^T x + b|}{\|w\|_2}$$

$\searrow$ distance of $x$ to the hyperplane $w^T x + b = 0$



$$w^T x = 0$$
$$w^T x + b = 0 \to$$
$$w$$
$$\frac{b}{\|w\|}$$
$$x$$
$$\frac{w \cdot x}{\|w\|}$$

$(\text{Assuming } w^T x > 0, b > 0)$

distance of $x$ to $w^T x + b = 0$ is

$$\frac{|w^T x + b|}{\|w\|}$$

SVM: Finds a hyperplane with the maximum margin.

For a given $h$ & dataset $S = \{x_1, \text{---} x_n\}$

$$\ell_h = \min_{i=1 \text{---} m} \ell_h(x_i)$$

The max-margin $\ell$ is

$$\rho = \max_{w,b \,:\, y_i(w^T x_i + b) \geq 0} \min_{i \in \{1 \dots m\}} \frac{|w^T x_i + b|}{\|w\|}$$

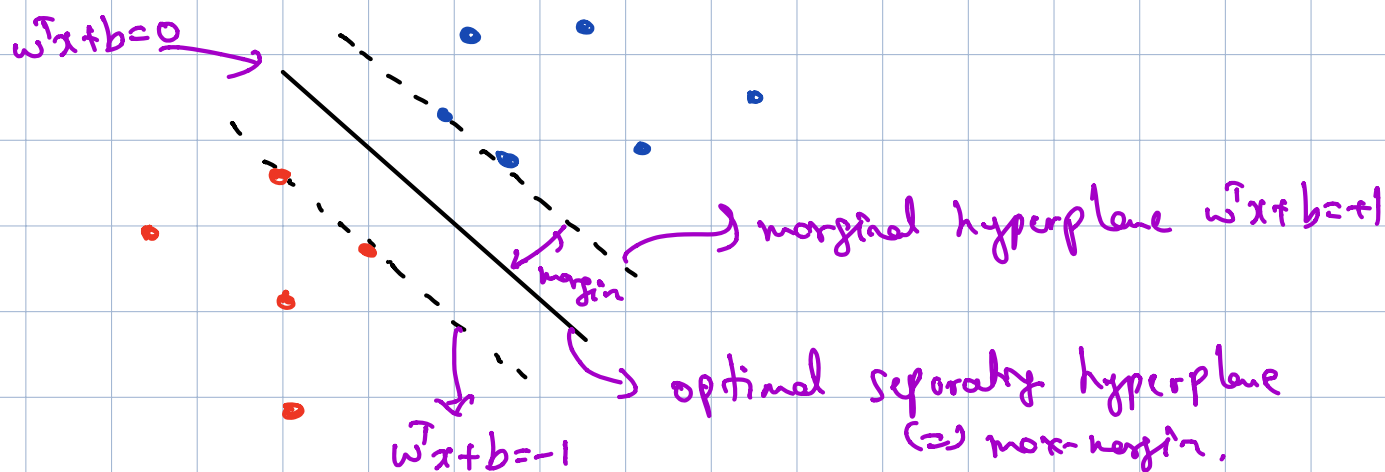$$\underset{*}{=} \max_{w,b} \min_{i \in \{1 \dots m\}} \frac{y_i(w^T x_i + b)}{\|w\|}$$

holds because the dataset is linearly separable, which implies $y_i(w^T x_i + b)$ is positive for the hyperplane defined by $(w, b)$ that achieves the maximum.

It is enough to consider scaled $(w, b)$ s.t.

$$\min_{i \in \{1 \dots m\}} y_i(w^T x_i + b) = 1$$

So, $\rho = \max_{\substack{w, b \,:\\ \min_{i \in \{1 \dots m\}} y_i(w^T x_i + b) = 1}} \frac{1}{\|w\|}$

$$= \max_{\substack{w, b \,:\\ \forall i, \, y_i(w^T x + b) \geq 1}} \frac{1}{\|w\|}$$



$w^T x + b = 0$

$\longrightarrow$ marginal hyperplane $w^T x + b = +1$

margin

$w^T x + b = -1$

$\longrightarrow$ optimal separating hyperplane
($=$) max-margin.

# Optimization problem underlying the SVM method:

$$\text{Maximizing} \quad \rho \propto \frac{1}{\|w\|} \quad (=) \quad \text{minimizing} \quad \frac{1}{2}\|w\|^2$$

$$\text{s.t.} \quad y_i(w^T x_i + b) \geq 1 \quad \forall i$$

## Primal problem!

$$\min_{w,b} \quad \frac{1}{2}\|w\|^2$$

$$\text{Subject to} \quad y_i(w^T x_i + b) \geq 1 \quad, i = 1 \cdots m$$

" Quadratic optimization problem with linear inequality constraints".

---

A brief tour of constrained optimization

Ref: Appendix B of FOML book

Constrained optimization problem!

$$\min_{x \in X} \quad f(x)$$
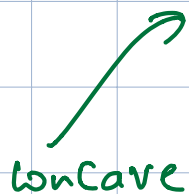
$$\text{s.t.} \quad g_i(x) \leq 0, \quad i = 1 \cdots m$$

Lagrangian $L(x, \alpha_1, --- \alpha_m)$ is defined by

$$L(x, \alpha_1, -- \alpha_m) = f(x) + \sum_{i=1}^{m} \alpha_i \, g_i(x)$$

Dual function:

$$F(\alpha_1, --- \alpha_m) = \inf_{x \in X} L(x, \alpha_1, --- \alpha_m), \text{ for any}$$
$$\alpha_1, -- \alpha_m \geqslant 0$$
$$= \inf_{x \in X} \left( f(x) + \sum_{i=1}^{m} \alpha_i \, g_i(x) \right)$$

concave

### Dual problem :-

$$\max_{\alpha_1, -- \alpha_m} F(\alpha_1, -- \alpha_m)$$
$$\text{s.t.} \quad \alpha_i \geqslant 0, \quad i = 1, -- m$$

Karash-Kuhn-Tucker (KKT) conditions :

Assume $f, g_i, i=1, --m$ are convex & differentiable

" constraint qualification" holds ( $\exists \, \bar{x} \in$ interior of $X$
s.t. $g_i(\bar{x}) < 0 \, \forall i$ )

$\bar{x}$ is a solution of the primal problem
if and only if

$$\exists (\bar{\lambda}_1, \ldots \bar{\lambda}_m), \quad \bar{\lambda}_i \geq 0, \forall i \quad s.t.$$

$$\nabla_x L(\bar{x}, \bar{\lambda}_1, \ldots \bar{\lambda}_m) = 0 \iff \nabla_x f(\bar{x}) + \sum_{i=1}^{m} \bar{\lambda}_i \nabla_x g_i(\bar{x}) = 0$$

$$g_i(\bar{x}) \leq 0, \quad i = 1 \ldots m \qquad \text{"feasibility"}$$

$$\bar{\lambda}_i \, g_i(\bar{x}) = 0, \quad i = 1 \ldots m \qquad \text{"complementary slackness"}$$

Note: ① Let $p^*$ be the optimal value of the primal problem & $d^*$ ⎽⎽⎽ " ⎽⎽⎽ " ⎽⎽⎽ dual problem

Weak duality: $\qquad d^* \leq p^*$

$$\left( p^* - d^* \right) \leftarrow \text{Duality gap}$$

Strong duality: $\qquad d^* = p^* \qquad (= \text{no duality gap})$

If the primal problem is a convex optimization problem
($=$ $f, g_i$ are convex) & a constraint qualification holds
($= \exists \bar{x} \in int(X)$ s.t $g_i(\bar{x}) < 0 \; \forall i$), then
" Strong duality" holds.

① If the constraints are linear, then the constraint qualification holds.

---

## Back to the SVM optimization problem in the linearly separable case:

Recall
$$\min_{w,b} \quad \frac{1}{2}\|w\|^2$$

$$\text{s.t.} \quad y_i(w^T x_i + b) \geq 1$$

Lagrangian:

$$L(w, b, \alpha_1, \ldots \alpha_m) = \frac{1}{2}\|w\|^2 - \sum_{i=1}^{m} \alpha_i\left(y_i(w^T x_i + b) - 1\right)$$

### KKT conditions:

(*) 

$$\nabla_w L = w - \sum_{i=1}^{m} \alpha_i y_i x_i = 0 \quad \Rightarrow \quad w = \sum_{i=1}^{m} \alpha_i y_i x_i$$

$$\nabla_b L = -\sum_{i=1}^{m} \alpha_i y_i = 0 \quad \Rightarrow \quad \sum_{i=1}^{m} \alpha_i y_i = 0$$

$$\forall i, \quad \alpha_i\left(y_i(w^T x_i + b) - 1\right) = 0 \quad \Rightarrow \quad \alpha_i = 0 \text{ (or)}$$
$$y_i(w^T x_i + b) = 1$$

Observe that

optimal $w$ is a linear combination of the data $x_1, --x_m$.

In fact, let $S = \{ i \mid \alpha_i \neq 0 \}$

Then, from $(*)$

$$w = \sum_{i \in S} \alpha_i y_i x_i$$

Vectors $x_i$ with $\alpha_i \neq 0$ are the support vectors.

From $(*)$, for the support vectors $y_i (w^T x_i + b) = 1$.

(or) the support vectors lie on the marginal hyperplanes.

Also, $b = y_i - w^T x_i$, where $x_i$ is a support vector

$$b = y_i - \sum_{j=1}^{n} \alpha_j y_j (x_j^T x_i)$$

**Dual optimization problem:**

$$F(\alpha_1, --\alpha_m) = \inf_{w, b} \left[ \frac{1}{2} \|w\|^2 - \sum_{i=1}^{m} \left( \alpha_i \left( y_i (w^T x_i + b) - 1 \right) \right) \right]$$

Note that if $\sum \alpha_i y_i \neq 0$, then $F(\alpha_1, -\alpha_m) = -\infty$.

So, lets impose $\sum \alpha_i y_i = 0$

The infimum is attained at $\quad w = \sum\limits_{i=1}^{m} \alpha_i y_i x_i$

$$F(\alpha_1, \ldots \alpha_m) = \frac{1}{2} \left\| \sum_{i=1}^{m} \alpha_i y_i x_i \right\|^2 - \sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j (x_i^T x_j)$$
$$- \sum_{i=1}^{m} \alpha_i y_i b + \sum_{i=1}^{\hat{m}} \alpha_i$$

$$F(\alpha_1, \ldots \alpha_m) = -\frac{1}{2} \sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j (x_i^T x_j) + \sum_{i=1}^{m} \alpha_i$$

Dual optimization problem:-

(**)$\left\{ \begin{array}{l} \underset{\alpha_1 - \alpha_m}{max} \quad -\frac{1}{2} \sum\limits_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j (x_i^T x_j) + \sum\limits_{i=1}^{m} \alpha_i \\[20pt] \text{S.t.} \quad \alpha_i \geq 0 \quad \text{and} \quad \sum\limits_{i=1}^{m} \alpha_i y_i = 0 \;, \; i=1 \ldots m \end{array} \right.$

(**) Is a constrained optimization problem, where the objective is quadratic and the constraints are linear.

The dimension of the dual problem is "m"

Can we QP solvers for (**).

"Strong duality" holds. Solving (★) gives $\alpha_1, \ldots, \alpha_m$, which can be used to determine the solution $w, b$ of the primal problem, leading to the following SVM classifier:

$$h(x) = \text{sign}(w^T x + b)$$

$$= \text{sign}\left(\sum_{i=1}^{m} \alpha_i y_i (x_i^T x) + b\right),$$

where $b = y_i - \sum_{j=1}^{m} \alpha_j y_j (x_j^T x_i)$.

---

Remark! " Margin $\ell$ of the SVM classifier ".

$$b = y_i - \sum_{j=1}^{m} \alpha_j y_j (x_j^T x_i).$$

$$\sum_{i=1}^{m} \alpha_i y_i \, b = \sum_{i=1}^{m} \alpha_i y_i^2 - \sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j (x_i^T x_j)$$
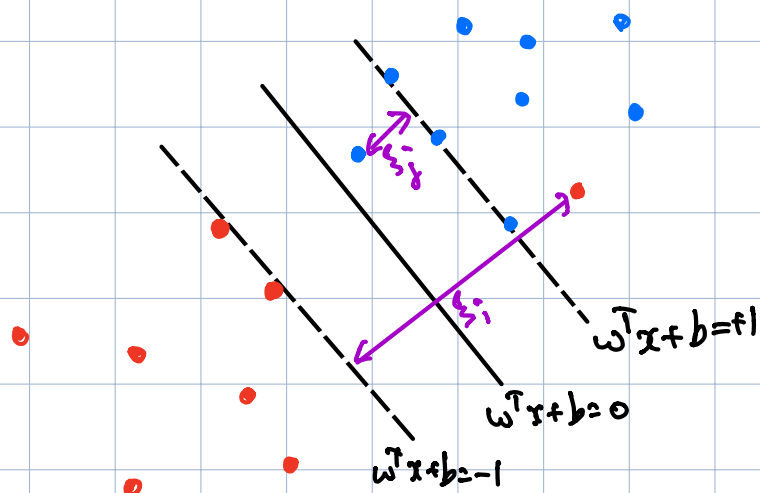
$$0 = \sum_{i=1}^{m} \alpha_i - \left(\sum_{i=1}^{m} \alpha_i y_i x_i\right)^T \left(\sum_{j=1}^{m} \alpha_j y_j x_j\right)$$

$$0 = \sum_{i=1}^{m} \alpha_i - \|w\|^2 \qquad - \text{①}$$

used ①
& $\alpha_i \geqslant 0$

$$\ell^2 = \frac{1}{\|w\|^2} = \frac{1}{\sum_{i=1}^{m} \alpha_i}$$

# SVM in the non-separable case (aka soft-margin SVM)



$$S = \{ x_1 -- x_m \}$$

Training data is not linearly separable, which is equivalent to

$\forall$ hyperplanes $w^T x + b = 0$, $\exists x_i \in S$ s.t.

$$y_i (w^T x_i + b) \not\geq 1. \quad ——(\ast)$$

A relaxed version of $(\ast)$ holds, i.e.,

$$\forall i = 1 -- m, \exists \xi_i \geq 0 \text{ such that}$$

$$y_i (w^T x_i + b) \geq 1 - \xi_i$$

$\xi_i \Leftarrow$ slack variables & measure by how much a point $x_i$ violates the separability constraint $y_i (w^T x + b) \geq 1$

If we ignore the outliers, then a margin of $\ell = \dfrac{1}{\|\omega\|}$ is achieved & hence, this margin is referred to as "soft margin".

## Primal optimization problem:

$$(***) \begin{cases} \underset{\omega, b, \xi}{\min} \quad \dfrac{1}{2}\|\omega\|^2 + C\sum_{i=1}^{m}\xi_i^p \\ \\ s.t. \quad y_i(\omega^T x_i + b) \geq 1 - \xi_i, \quad i=1\dots m \\ \\ \quad \xi_i \geq 0, \quad i=1\dots m \end{cases}$$

Objective here is $\min \dfrac{1}{2}\|\omega\|^2 + C\sum_{i=1}^{m}\xi_i^p$ ,

$\min \|\omega\|^2$ $(=)$ maximizing the margin

minimize the slack. (slack is due to outliers)

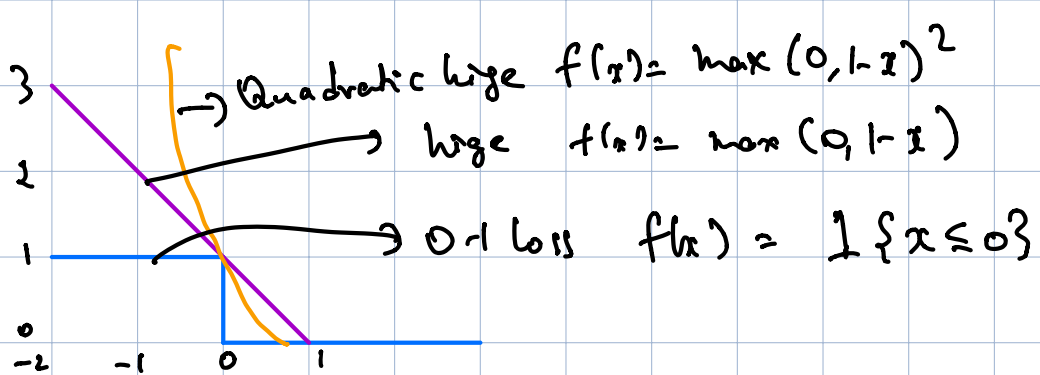where $p \geq 1$

$(***)$ is a convex optimization problem

since $\xi \rightarrow \sum_{i=1}^{m}\xi_i^p = \|\xi\|_p^p \quad \xi = (\xi_1 \dots \xi_m)$

is a convex function

Choice of p:            $p=1$       and       $p=2$

Hinge loss                      Quadratic hinge loss

→ Quadratic hinge $f(x) = \max(0, 1-x)^2$

→ hinge $f(x) = \max(0, 1-x)$

→ 0-1 loss $f(x) = \mathbb{1}\{x \leq 0\}$

Hinge loss is popular.
We focus on $p=1$.

→ allows trade off between max-margin & reducing slack

$(***) \to$
$$\begin{cases} \min_{w, b, \xi} & \frac{1}{2}\|w\|^2 + C\sum_{i=1}^{m}\xi_i \\ \text{s.t.} & y_i(w^T x_i + b) \geq 1 - \xi_i, \\ & \xi_i \geq 0 \end{cases} \quad i=1 \cdots m$$

Recipe for solving $(***)$ is

① Apply KKT conditions

② Work out dual function

③ formulate + Solve dual optimization problem

④ We soln of dual problem to obtain the soft-margin SVM hyperplane classifier.

Lagrangian
$$L(w, b, \xi, \alpha, \beta) = \frac{1}{2}\|w\|^2 + C\sum_{i=1}^{m}\xi_i - \sum_{i=1}^{m}\alpha_i\left(y_i(w^T x_i + b) - 1 + \xi_i\right)$$
$$- \sum_{i=1}^{m}\beta_i\xi_i$$

We apply KKT conditions

$$\nabla_w L = w - \sum_{i=1}^{m} \alpha_i y_i x_i = 0 \quad \Rightarrow \quad w = \sum_{i=1}^{m} \alpha_i y_i x_i$$

$$\nabla_b L = -\sum_{i=1}^{m} \alpha_i y_i = 0 \quad \Rightarrow \quad \sum_{i=1}^{m} \alpha_i y_i = 0$$

$$\nabla_{\xi_i} L = C - \alpha_i - \beta_i = 0 \quad \Rightarrow \quad \alpha_i + \beta_i = C$$

$$\forall i, \quad \alpha_i \left( y_i (w^T x_i + b) - 1 + \xi_i \right) = 0 \quad \Rightarrow \quad \alpha_i = 0 \quad (or)$$

$$y_i (w^T x_i + b) = 1 - \xi_i$$

$$\forall i, \quad \beta_i \xi_i = 0 \quad \Rightarrow \quad \beta_i = 0 \quad (or) \quad \xi_i = 0$$

**Remarks:**

① Optimal $w$ has the same expression as before
(= lin. sep. case)

② $$w = \sum_{i=1}^{m} \alpha_i y_i x_i$$

$x_i$ appears only if $\alpha_i \neq 0$

Such $x_i$'s are the support vectors.

for a support vector $x_i$, we have
$$y_i (w^T x_i + b) = 1 - \xi_i$$

If $\xi_i = 0$, then $y_i (w^T x_i + b) = 1$ & $x_i$ is on one of the
the marginal hyperplanes ($w^T x + b = \pm 1$)

If $\xi_i \neq 0$, then $x_i$ is an outlier and
$\beta_i = 0$ which implies $\alpha_i = C$.

## Dual optimization problem:

$$F = \frac{1}{2} \left\| \sum_{i=1}^{m} \alpha_i y_i x_i \right\|^2 - \sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j (x_i^T x_j)$$
$$- \sum_{i=1}^{m} \alpha_i y_i b + \sum_{i=1}^{m} \alpha_i$$

$$F = \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j (x_i^T x_j)$$

Dual function "is the same as" before $\left( = \substack{\text{lin. sep.} \\ \text{case}} \right)$

However, in addition to $\alpha_i \geq 0$ & $\sum \alpha_i y_i = 0$,

we require $\alpha_i \leq C$.

So, the dual optimization problem is

(\*\*\*\*) $\begin{cases} \max\limits_{\alpha_1 \cdots \alpha_m} \sum\limits_{i=1}^{m} \alpha_i - \frac{1}{2} \sum\limits_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j (x_i^T x_j) \\ \\ \text{s.t.} \quad \alpha_i \geq 0, \quad \alpha_i \leq C, \quad i = 1 \cdots m \\ \qquad\quad \sum\limits_{i=1}^{m} \alpha_i y_i = 0 \end{cases}$

(\*\*\*\*) is similar to (\*\*) (= dual problem in the separable case), with an additional constraint $\alpha_i \leq C$.

The complexity of $^{solve}$ (✱✱✱✱) is comparable to that of (✱✱); & one could use a QP solver in either case.

The solution $(\alpha_1 , -- \alpha_m)$ of (✱✱✱✱) can be used to define the soft margin SVM hyperplane as follows:

$$h(x) = \text{sign}(\vec{w}^T x + b)$$

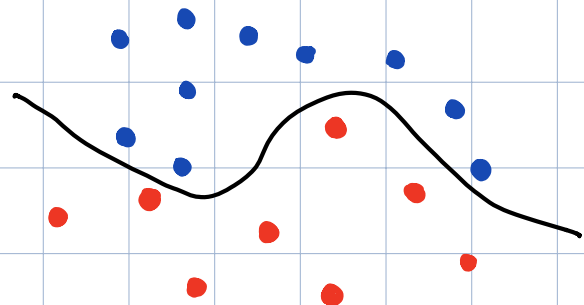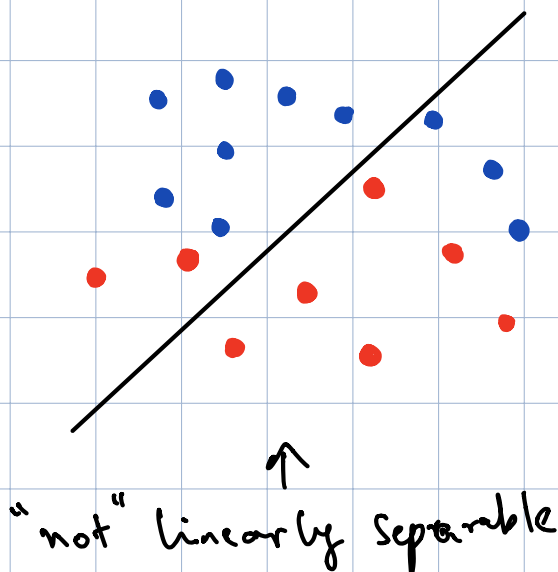$$= \text{sign}\left( \sum_{i=1}^{m} \alpha_i y_i (x_i^T x) + b \right),$$

where

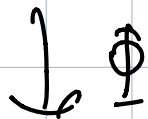$$b = y_i - \sum_{j=1}^{m} \alpha_j y_j (x_j^T x_i),$$

for any $x_i$ with $0 < \alpha_i < C$.

---

Kernel methods

↑
"not" linearly separable

↑
non-linear decision boundary

Idea:   Take data $(x_i)$

$$\downarrow \Phi$$

High-dimensional space $H$

$\searrow$ run SVMs in $H$, since the data could be linearly separable in $H$.

SVMs in high-dimensions:

$$h(x) = \text{sign}(w^T x + b)$$

$$= \text{sign}\left( \sum_{i=1}^{m} \alpha_i y_i \, \phi(x_i)^T \phi(x) + b \right),$$

is a high-dimensional feature mapping $\phi(x) \in H$, high-dimensional space.

where

$$b = y_i - \sum_{j=1}^{m} \alpha_j y_j \left( \phi(x_j)^T \phi(x_i) \right)$$

for any $x_i$ with $0 < \alpha_i < C$.

Observe! feature $x_i$ appear only through inner products.

Using the kernel trick, the inner products can be computed efficiently even in high dimensions.

Def (Kernel)

$\rightarrow$ input space (say some subset of $\mathbb{R}^d$) e.g.

A function $K: X \times X \rightarrow \mathbb{R}$ such that

$$\forall x, x' \in X, \quad K(x, x') = \langle \phi(x), \phi(x') \rangle$$

for some function $\phi: X \rightarrow H$.

H → Hilbert space

( Hilbert space H is a vector space that is equipped
with an inner product & is complete
(= all Cauchy sequences converge).

The norm induced by the inner product is

$$\|x\|_H = \sqrt{\langle x, x \rangle} \quad \forall x \in H \quad )$$

**Note:** There are kernels, where $K(x, x')$ can be
computed in $O(d)$ ($d$ = input feature dimension),
while computing $\langle \phi(x), \phi(x') \rangle$ is
$O(\dim(H))$, with $\dim(H) >> d$

**Question:** Given a function $K$, can we infer that

$K$ is a kernel, (or) equivalently there is a
Hilbert space $H$ & feature mapping $\phi$ s.t.

$(*) \longrightarrow K(x, x') = \langle \phi(x), \phi(x') \rangle$ ?

Yes, if $K$ is "positive definite symmetric" (PDS)

If K is PDS, then $\exists\ \mathcal{H}\ \&\ \phi$ s.t.
(*) holds.

So, one need not ~~"define"~~ or "compute" $\phi$.

K being PDS assures existence of a $\phi$.

## Definition (PDS)

A kernel $K: X \times X \to \mathbb{R}$ is PDS if for any $\{x_1, \ldots x_m\} \subseteq X$, the matrix

$$K = \left(\left[K(x_i, x_j)\right]\right)_{i,j=1\sim m} \quad \text{is Symmetric positive}$$

Semi-definite (SPSD)

kernel
matrix

K is SPSD if one of the following conditions holds:

① Eigenvalues of K are non-negative

② for any vector $c = (c_1, \ldots c_m)^T$,
$$c^T K c = \sum_{i,j} c_i c_j K(x_i, x_j) \geq 0$$

# Examples of PDS Kernels

① Polynomial kernel

Fix a constant $c > 0$. A polynomial kernel of degree $\beta$ is given by

$$K(x, x') = (x^T x' + c)^\beta, \quad \forall x, x' \in \mathbb{R}^d$$

Special case with $\beta = 2$, $d = 2$, $\forall x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, x' = \begin{bmatrix} x'_1 \\ x'_2 \end{bmatrix}$
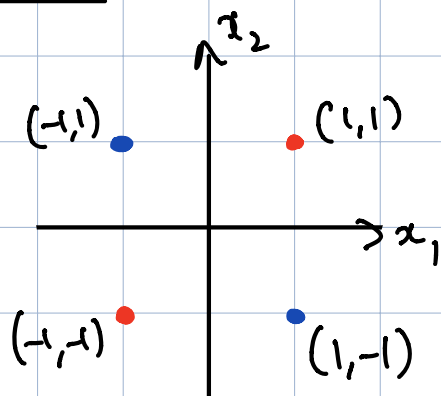
$$K(x, x') = (x^T x' + c)^2$$

$$= (x_1 x'_1 + x_2 x'_2 + c)^2$$

$$= \begin{bmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2} x_1 x_2 \\ \sqrt{2c} x_1 \\ \sqrt{2c} x_2 \\ c \end{bmatrix}^T \begin{bmatrix} x_1'^2 \\ x_2'^2 \\ \sqrt{2} x'_1 x'_2 \\ \sqrt{2c} x'_1 \\ \sqrt{2c} x'_2 \\ c \end{bmatrix}$$

$$= \phi(x)^T \phi(x')$$

# Achieving linear separation using a polynomial Kernel:



XOR problem
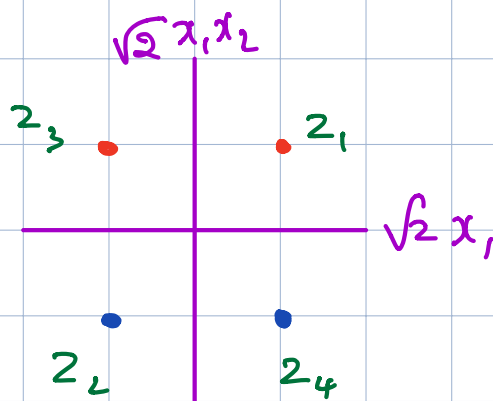"not" linearly separable

Use a polynomial kernel with $d=2$, $C=1$

$$(1,1) \to (1, 1, \sqrt{2}, \sqrt{2}, \sqrt{2}, 1) \, z_1$$

$$(-1,1) \to (1, 1, \sqrt{2}, -\sqrt{2}, -\sqrt{2}, 1) \, z_2$$

$$(-1,-1) \to (1, 1, -\sqrt{2}, -\sqrt{2}, \sqrt{2}, 1) \, z_3$$

$$(1,-1) \to (1, 1, -\sqrt{2}, \sqrt{2}, -\sqrt{2}, 1) \, z_4$$

Linear Separation $\longrightarrow$
using a polynomial
Kernel



Note: Polynomial kernel is PDS since we wrote as a inner product with an explicit $\phi$.

Example 2: Gaussian Kernels

$$\forall x, x' \in \mathbb{R}^d, \quad K(x, x') = \exp\left( -\frac{\|x - x'\|^2}{2\sigma^2} \right)$$

Gaussian Kernel is PDS. (Can be shown using the normalization property of PDS kernels)

**Example 3:** Sigmoid kernels

$$\forall x, x' \in \mathbb{R}^d, \quad K(x, x') = \tanh\left(a(x^T x') + b\right)$$

note: sigmoid kernels + SVM = simple neural network

---

**Claim without proof:**

Let $K: X \times X \to \mathbb{R}$ be a PDS kernel.
Then, there exists a Hilbert space $H$ and a
  mapping $\phi: X \to H$ s.t.

$$\forall x, x' \in X, \quad K(x, x') = \langle \phi(x), \phi(x') \rangle.$$

---

Verifying that the Gaussian kernel is PDS.

Normalized kernel $K'$ associated with a
  PDS kernel $K$ is

$$\forall x, x' \in X, \quad K'(x, x') = \begin{cases} 0 & \text{if } K(x,x) = 0 \text{ (or) } K(x',x') = 0 \\ \dfrac{K(x,x')}{\sqrt{K(x,x) \, K(x',x')}} & \text{else} \end{cases}$$

By definition, $K'(x, x) = 1 \quad \forall x \in X.$

A Gaussian Kernel can be seen as the normalized kernel of the kernel $K'(x,x') = \exp\left(\frac{x^T x'}{\sigma^2}\right)$

This can be argued as follows:

$$\frac{K'(x,x')}{\sqrt{K'(x,x) \, K'(x',x')}} = \frac{\exp\left(\frac{x^T x'}{\sigma^2}\right)}{\exp\left(\frac{\|x\|^2}{2\sigma^2}\right) \exp\left(\frac{\|x'\|^2}{2\sigma^2}\right)}$$

$$= \exp\left(-\frac{\|x-x'\|^2}{2\sigma^2}\right)$$

in the Gaussian Kernel (obtained by normalizing $K'$)

<u>Claim without proof</u>: If $K'$ is PDS, then its normalized variant is PDS as well.

To see that $K'(x,x') = \exp\left(\frac{x^T x'}{\sigma^2}\right)$ is PDS,

$$K'(x,x') = \sum_{j=0}^{\infty} \frac{(x^T x')^j}{\sigma^{2j} \; j!} \longrightarrow \text{polynomial kernel}$$

$K'$ is a positive linear combination of polynomial kernels.

**Claim:** If $K, K'$ are PDS kernels, then $K + K'$ is a PDS kernel.

**Why?**
$$c \in \mathbb{R}^{m \times 1}$$

$$c^T K c \geq 0, \quad c^T K' c \geq 0 \quad (\text{since } K, K' \text{ are PDS})$$

$$\Rightarrow \quad c^T (K + K') c \geq 0$$

**Claim:** PDS kernels closed under pointwise limit & composition with $f: x \to \sum_{j=0}^{\infty} a_j x^j$, $a_j \geq 0$. $\longleftarrow$ check this

Hence, Gaussian kernel is PDS.

---

## SVMs with PDS Kernels:

Idea: Replace $x^T x'$ with $K(x, x')$

Dual optimization problem after incorporating kernel is

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

S.t. $\quad 0 \leq \alpha_i \leq C, \quad i = 1 \cdots m$

$$\sum_{i=1}^{m} \alpha_i y_i = 0$$

Solving the optimization problem would lead to the following classifier:

$$h(x) = \text{sign}\left( \sum_{i=1}^{m} \alpha_i y_i k(x_i, x) + b \right),$$

where 
$$b = y_i - \sum_{j=1}^{m} \alpha_j y_j k(x_j, x_i)$$

$$\text{for} \quad x_i \text{ s.t. } 0 < \alpha_i < C$$

---

## Kernel Ridge Regression

$$\left[ \text{Sec 11.3 of FOML book} \right]$$

Input space $\mathcal{X} \subseteq \mathbb{R}^d$
Feature mapping $\phi : \mathcal{X} \to \mathbb{R}^{d'}$

Linear hypothesis set:

$$\left\{ h \mid h(x) = w^T \phi(x), \; w \in \mathbb{R}^{d'} \right\}$$

<u>linear regression (recall)</u>  $S = \{ (x_i, y_i), i = 1 \ldots m \}$

$$\min_{w} \frac{1}{m} \sum_{i=1}^{m} (w^T \phi(x_i) - y_i)^2$$

$$= \min_{w} \left\{ J(w) := \frac{1}{m} \| Aw - Y \|^2 \right\},$$

where   A is the feature matrix with rows $\phi(x_i)^T$,

$Y$ is a vector with components $y_i$

$J(w)$ is minimized by the solution to

$$A^T A W = A^T Y$$

---

## Ridge regression & its connection to kernels :-

$$\min_{w} \sum_{i=1}^{m} (w^T \phi(x_i) - y_i)^2 + \lambda \|w\|^2$$

$$= \min_{w} \quad \underbrace{\|A W - Y\|^2 + \lambda \|w\|^2}_{\overset{..}{J}(w)}$$

$$\nabla J(w) = 0$$

$$\Rightarrow \quad (A^T A + \lambda I) W = A^T Y$$

$$\text{(or)} \qquad W = (A^T A + \lambda I)^{-1} A^T Y \quad \underline{\qquad} \quad (1)$$

$$\underset{\substack{\downarrow \\ \text{is invertible because}}}{}$$

$$\qquad\qquad A^T A \text{ is positive semidefinite}$$
$$\& \lambda > 0$$

An equivalent formulation for ridge regression:

$$\min_{w} \sum_{i=1}^{m} (w^T \phi(x_i) - y_i)^2 \qquad \left.\begin{array}{c} \\ \\ \end{array}\right\} \rightarrow (01)$$

$$\text{subject to} \quad \|w\|^2 \le \tau^2$$

(01) In a constrained optimization problem
with convex objective as well as convex constraints
(Why?)

(01) can be re-written as

$$\min_{w} \quad \sum_{i=1}^{m} \xi_i^2$$

$$\text{Subject to} \quad \|w\|^2 \leq \Lambda^2, \text{ and}$$
$$\xi_i = y_i - w^T \phi(x_i), \quad i=1 \dots m$$

$\left.\right\} - (02)$

(02) $\Rightarrow$ Convex optimization problem

Lagrangian

$L(\xi, w, \alpha', \lambda)$

$\boxed{\begin{array}{l} \xi = (\xi_1, ---\xi_m)^T \\ \alpha' = (\alpha'_1, ---\alpha'_m)^T \end{array}}$

$$= \sum_{i=1}^{m} \xi_i^2 + \sum_{i=1}^{m} \alpha'_i \left(y_i - \xi_i - w^T \phi(x_i)\right) + \lambda\left(\|w\|^2 - \Lambda^2\right)$$

$\longrightarrow$ _Lagrange multipliers_

Applying KKT conditions, we obtain

$$\nabla_w L = -\sum_{i=1}^{m} \alpha'_i \phi(x_i) + 2\lambda w = 0 \quad \Rightarrow \quad w = \frac{1}{2\lambda} \sum_{i=1}^{m} \alpha'_i \phi(x_i)$$

$$\nabla_{\xi_i} L = 2\xi_i - \alpha_i' = 0 \quad \Rightarrow \quad \xi_i = \frac{\alpha_i'}{2}$$

$$\alpha_i' \left( y_i - \xi_i - w^T \phi(x_i) \right) = 0, \quad i = 1 \cdots m$$

$$\lambda \left( \|w\|^2 - \Lambda^2 \right) = 0$$

---

Plugging in expressions for $w$ & $\xi_i$ from KKT conditions
into the Lagrangian

$$\sum_{i=1}^{m} \frac{\alpha_i'^2}{4} + \sum_{i=1}^{m} \alpha_i y_i - \sum_{i=1}^{m} \frac{\alpha_i'^2}{2} - \frac{1}{2\lambda} \sum_{i,j=1}^{m} \alpha_i' \alpha_j' \phi(x_i)^T \phi(x_j)$$

$$+ \lambda \left( \frac{1}{4\lambda^2} \left\| \sum_{i=1}^{m} \alpha_i' \phi(x_i) \right\|^2 - \Lambda^2 \right)$$

$$= -\frac{1}{4} \sum_{i=1}^{m} \alpha_i'^2 + \sum_{i=1}^{m} \alpha_i' y_i - \frac{1}{4\lambda} \sum_{i,j=1}^{m} \alpha_i' \alpha_j' \phi(x_i)^T \phi(x_j)$$

$$- \lambda \Lambda^2$$

Make the substitution $\alpha_i' = 2\lambda \alpha_i$

$$= -\lambda^2 \sum_{i=1}^{m} \alpha_i^2 + 2\lambda \sum_{i=1}^{m} \alpha_i y_i - \lambda \sum_{i,j=1}^{m} \alpha_i \alpha_j \phi(x_i)^T \phi(x_j)$$

$$- \lambda \Lambda^2$$

The dual optimization problem is

$$\underset{\alpha \in \mathbb{R}^m}{\text{argmax}} \left( -\lambda^2 \sum_{i=1}^m \alpha_i^2 + 2\lambda \sum_{i=1}^m \alpha_i y_i - \lambda \sum_{i,j=1}^m \alpha_i \alpha_j \phi(x_i)^T \phi(x_j) \right.$$
$$\left. - \lambda \Lambda^2 \right)$$

$$= \underset{\alpha \in \mathbb{R}^m}{\text{arg max}} \left( -\lambda \sum_{i=1}^m \alpha_i^2 + 2 \sum \alpha_i y_i - \underbrace{\sum_{i,j=1}^m \alpha_i \alpha_j \phi(x_i)^T \phi(x_j)}_{} \right)$$
$$= G(\alpha)$$

$$\underset{\alpha}{\max} \quad G(\alpha)$$

$$= \underset{\alpha}{\max} \quad -\lambda \alpha^T \alpha + 2 \alpha^T Y - \alpha^T (A A^T) \alpha$$

$$= \underset{\alpha}{\max} \quad -\alpha^T (K + \lambda I) \alpha + 2 \alpha^T Y$$

<span style="color:green">↓</span>

<span style="color:green">$K = A A^T$ is the kernel matrix</span>

Optimizing the dual

$$\nabla G(\alpha) = 0$$

$$(=) \quad 2(K + \lambda I)\alpha = 2Y$$

$$\alpha = (K + \lambda I)^{-1} Y \qquad \sim (2)$$
$$\longrightarrow \text{is invertible (why?)}$$

Using KKT conditions,

$$w = \sum_{i=1}^{m} \alpha_i \phi(x_i) = A^T \alpha = A^T (K + \lambda I)^{-1} y$$

Linear hypothesis

$$h(x) = w^T \phi(x)$$

$$h(x) = \sum_{i=1}^{m} \alpha_i K(x_i, x)$$

Any PDS Kernel can be used to arrive at this predictor.

| | Solving primal i.e., (1) | Solving dual, i.e., (2) |
|---|---|---|
| Computational Cost | $A^T A : O(m d'^2)$ <br><br> $(A^T A + \lambda I)^{-1}:$ <br> $O(d'^3)$ <br><br> multiply with $A^T$: <br> $O(d')$ | Let $\kappa$ be the $\widetilde{\text{max}}$ cost of computing $K(x, x') \ \forall x, x'$ <br><br> Kernel matrix is computed in $O(\kappa m^2)$ <br><br> Inverting $(K + \lambda I):$ <br> $O(m^3)$ <br><br> multiplication with $y:$ <br> $O(m^2)$ |
| Total cost: | $O(m d'^2 + d'^3)$ | Total cost: $O(\kappa m^2 + m^3)$ |

$$O((d')^3) \quad vs \quad O(m^3)$$

If $\phi$ is a mapping onto a high dimensional feature space & if the # of training samples is moderate, $d' >> m$, then solving the dual is computationally advantageous.

Prediction cost! $w^T \phi(x)$ computed in $O(d')$ for the primal

In case of the dual,

computing $(K(x_1, x), \, - \, - \, - \, - \, , \, K(x_m, x)) = \psi$

for a given $x$ is $O(km)$

& $\psi^T \alpha$ is $O(m)$

So, total prediction cost $O(km)$