→ Convex Optimization

→ SVM , KKT conditions to SVM
   (derivation via margin maximization)

## Variants of SVM

$$\min_{\substack{w, b \in \mathbb{R} \\ \in \mathbb{R}^d}} \frac{\|w\|^2}{2}$$
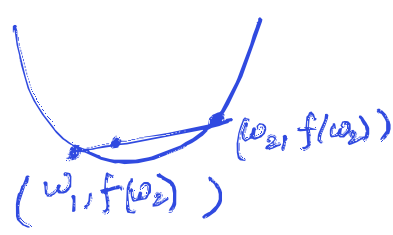
$$s.t. \quad y_i(w \cdot x_i + b) \geq 1$$

$d = dim(x)$

Margin $\dfrac{1}{\|w\|}$

# Convex optimization

Convex function $f$  a convex set $\omega$
for any  $\omega_1, \omega_2 \in \omega$,
$$f(\alpha \omega_1 + (1-\alpha) \omega_2) \leq \alpha f(\omega_1)$$
$$+ (1-\alpha) f(\omega_2)$$

$< $ strictly convex


$(\omega_2, f(\omega_2))$
$(\omega_1, f(\omega_2))$

**Thm:**  Suppose a convex function $f: \omega \to \mathbb{R}$ has a minimum. Then, the set of minima
forms a convex subset of $\omega$.
If $f$ is strictly convex, there is a
unique minimum.

## Proof:
Let $M_f \subseteq \omega$ be the set of minima of
$f$. Let $m_f$ be the minimum value of $f$ on $\omega$.
Let  $\omega, \omega' \in M_f$ , for any $\alpha \in [0,1]$,
$\alpha \omega + (1-\alpha) \omega' \in M_f$.

$$f(\alpha \omega + (1-\alpha)\omega') \leq \alpha f(\omega) + (1-\alpha)f(\omega')$$
$$= \alpha m_f + (1-\alpha) m_f$$
$$= m_f$$

$$\Rightarrow \quad \alpha \omega + (1-\alpha)\omega' \in M_f$$

when $f$ is strictly convex, $\omega = \omega'$.
$\qquad M_f$ is a singleton set.


# KKT theorem for constrained optimization

$$\min_{\omega \in \mathbb{R}^d} f(\omega) = P_{opt}$$
$$\text{s.t.} \quad g_i(\omega) \leq 0 \quad i = [m]$$  primal

Lagrangian: $\mathcal{L}(\omega, \lambda) = f(\omega) + \sum_{i=1}^{m} \lambda_i g_i(\omega)$

dual

$$\lambda \in \mathbb{R}^m \qquad \lambda = [\lambda_1, \ldots, \lambda_m]^T \in \mathbb{R}^m$$

$f(\omega)$   $F(\lambda)$



Dual problem:
$$\max_{\lambda \in \mathbb{R}^m} F(\lambda) := \inf_{\omega} \mathcal{L}(\omega, \lambda) = d_{opt}$$
$$\text{s.t.} \quad \lambda_i \geq 0 \qquad i \in [m].$$

Weak duality: $d_{opt} \leq P_{opt}$
with equality for convex problem.

Assumption (1) $f$ is convex, $g$ is convex & diff.
KKT:  There exists a $\lambda \in \mathbb{R}^m$ s.t. at
a minimum $\omega \in \mathbb{R}^d$ of $f$, the following
conditions are satisfied:

(i)  $\nabla_\omega \mathcal{L}(\omega, \lambda) = 0$

(ii)  $\nabla_\lambda \mathcal{L}(\omega, \lambda) = \left[ g_i(\omega) \leq 0 \right]_{i=1,\ldots,m}$

(iii) Complementarity constraints:
$\qquad$ either $\lambda_i = 0$ or $g_i(\omega) = 0$.
$\qquad\qquad$ for every $i = 1, \ldots, m$.

(2) Constraints are qualified: $\exists \ \omega \in int(\omega)$
$\text{s.t.} \quad g_i(\omega) < 0 \quad \forall \ i \in [m]$.

or  $g_i(\omega) \leq 0$ and $g_i$ is affine

(Slater's conditions)  $\qquad g_i(\omega) = x \cdot \omega$
$\qquad\qquad\qquad\qquad\qquad + b$

Ref: Chapter 6 of "Learning with Kernels"
$\qquad$ Smola & Scholkopf.

# SVM

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

$$f(\omega) = \frac{\|\omega\|^2}{2}$$

$$\text{s.t.} \quad y_i(\omega \cdot x_i + b) \geq 1.$$
$$i \in [m]$$

$$g_i(\omega) = 1 - y_i(\omega \cdot x_i + b)$$

$$\mathcal{L}(\omega, b, \lambda) = \frac{\|\omega\|^2}{2} - \sum_{i=1}^{m} \lambda_i (y_i(\omega \cdot x_i + b) - 1)$$

At $\min \omega, b, \; \exists \lambda \geq 0$

$$\left. \begin{array}{l} \nabla_\omega \mathcal{L}(\omega, b, \lambda) = 0 \\ \nabla_b \mathcal{L}(\omega, b, \lambda) = 0 \end{array} \right]$$

$$\omega = \sum_{i=1}^{m} \lambda_i y_i x_i$$

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

Comp. cond:
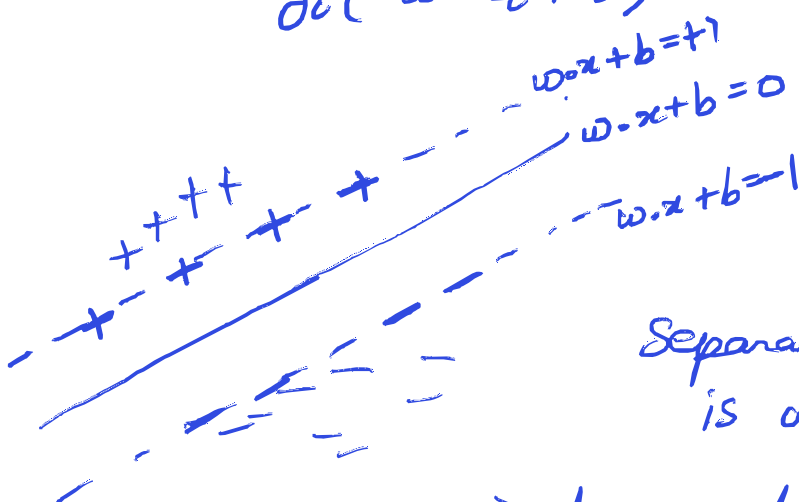$$\lambda_i = 0 \quad \text{or} \quad y_i(\omega \cdot x_i + b) = 1$$
$$i \in [m].$$

Last class: Force + torque balance

$x_i$ for non zero $\lambda_i \neq 0$ are called support vectors.

For support vectors $x_i$,
$$y_i(\omega \cdot x_i + b) = 1$$



$\omega \cdot x + b = +1$

$\omega \cdot x + b = 0$

$\omega \cdot x + b = -1$

Separating hyperplane is unique.

$> d$ support vectors, $\omega$ is not unique.

$$\omega = \sum_{i=1}^{m} \lambda_i y_i x_i \qquad \sum_{i=1}^{m} \lambda_i y_i = 0.$$

$$\lambda_i \neq 0, \qquad y_i(\omega \cdot x_i + b) = 1$$

**Dual problem**

$$\inf_{\omega, b} \mathcal{L}(\omega, b, \lambda) = \inf_{\omega, b} \frac{\|\omega\|^2}{2} - \sum_{i=1}^{m} \lambda_i \underbrace{(\omega \cdot x_i + b)y_i}_{-1}$$

$$\max_{\lambda} F(\lambda) = \frac{1}{2} \sum_{i,j=1}^{m} \lambda_i \lambda_j y_i y_j \, x_i \cdot x_j -$$
$$\sum_{i=1}^{m} \underline{\lambda_i \left( y_i \left( \sum_{j=1}^{m} \lambda_j y_j x_j \cdot x_i \right. \right.}$$
$$\left. \left. + b \right) - 1 \right)$$

$$\lambda_i \geqslant 0 \qquad i \in [m]$$

$$\max_{\lambda} \quad -\frac{1}{2} \sum_{i,j=1}^{m} \lambda_i \lambda_j y_i y_j \, x_i \cdot x_j + \sum_{i=1}^{m} \lambda_i$$
$$\text{s.t.} \quad \lambda_i \geqslant 0 \quad i \in [m].$$

Ex: Convex problem
Remark:
To evaluate objective function, we only
need dot products on data space

• $h_{SVM}(z) = \text{sgn}(\omega \cdot z + b)$
$$= \text{sgn}\left( \sum_{i=1}^{m} \lambda_i y_i x_i \cdot z + b \right)$$

output of SVC also needs only
dot products

$\Rightarrow$ So can replace $x_i \cdot z \longrightarrow \varkappa(x_i, z)$
to get nonlinear classifiers.
$$h_{SVM}(x) = \text{sgn}\left( \sum_{i=1}^{m} \lambda_i y_i \, \varkappa(x_i, z) + b \right)$$

When
$\varkappa$ PD kernel $\Rightarrow$ linear classifier
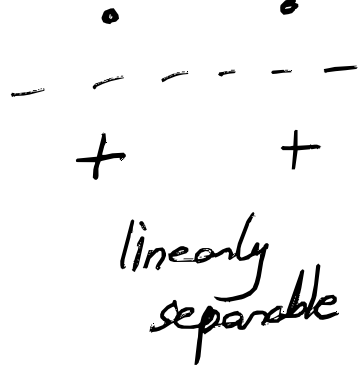with features $\Phi(x)$ s.t.
$$\Phi(x_i) \cdot \Phi(x_j) = \varkappa(x_i, x_j)$$

($\ell^2$ inner product if $\Phi(x)$ are
inf. dimensional).

$$x \qquad \longrightarrow \qquad \Phi(x)$$
Data                      feature
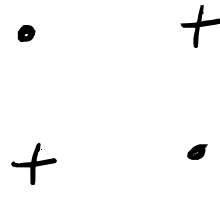(could be
finite- but
high-dimensional
OR
inf-dimensional)

SVM $\longrightarrow$ Kernel SVM

linear                nonlinear classifier
classifier               on $\mathcal{X}$
on $\mathcal{X}$
                      linear classifier
                      on feature space.

e.g.
XOR function



lineanly                non-linearly
separable               separable

$$\mathbb{R}^2 \longrightarrow \mathbb{R}^6$$

$$\begin{bmatrix} x^{(1)} \\ x^{(2)} \end{bmatrix} \quad \Rightarrow \quad \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ x^{(1)^2} \\ x^{(2)^2} \\ x^{(1)} x^{(2)} \\ c \end{bmatrix}$$

$$\varkappa(x_i, x_j) = (x_i \cdot x_j + c)^2$$
$$\text{polynomial kernel}$$

Other kernels        hyper parameter

$$\varkappa(\vec{x}, y) = e^{-\frac{\|x-y\|^2}{2q^2}} \qquad \varkappa$$

Gaussian   $\|x-y\|^2$ for given data points

$$O(\|x-y\|^2) = O(q^2)$$
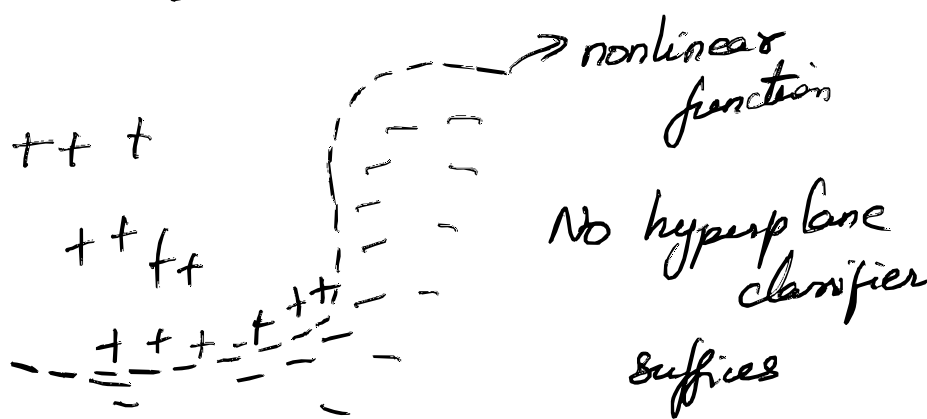
big O
( inf. dimensional feature space)

- Sigmoid kernel    (Neural network)

$$\varkappa(x, y) = \tanh(\sigma(x \cdot y) + b)$$
$$\sigma, b$$

### Beyond linear classifiers



$\rightarrow$ nonlinear function

No hyperplane classifier suffices

Kernel classifiers are still linear classifiers in feature space

Cover's **theorem**:
   m points in   d dimensions in "general position"

how many "labelings" / sets of m points are linearly separable?

# Subsets of m points / in dimensions that are linearly separable = C(m, d)
           ↑ labeling

- if $m \leq d+1$, $c(m, d) = 2^m$

- if $m > d+1$,
$$C(m, d) = 2 \sum_{i=0}^{d} \binom{m-1}{i}$$
         Sub exponential in m.

| Example | 2D | does lin. class. exist |
|---|---|---|



A    B   C

$+1 +1 +1$
$-1 -1 -1$ } Yes

$+1 +1 -1$    Yes
$\binom{3}{2}$

$+i \quad +1$

$\cdots -1$

$+1 \ -1 \ -1$    Yes
$\binom{3}{2}$

$$8 = 2^3$$



$+1$
$-1$    $+1$
A    $-1$
      B

$+1$    $+1$
D     C
$-1$    $-1$

$2^4 = 16$ subsets or labelings

\# pts $\rightarrow$ dim
$C(\overset{\downarrow}{4}, 2)$

$$= C(3, 2) + C(4, 1)$$
     $\|$        $\|$
     $8$        $6$

$$= 14$$

$m = 4$
$d = 2$

$$2 \sum_{i=0}^{d} \binom{m-1}{i} = 2\left(\binom{3}{0} + \binom{3}{1} + \binom{3}{2}\right)$$

$$= 2(1 + 3 + 3) = 14$$

Ex:
$$C(m, d) = c(m-1, d) + C(m, d-1)$$

Fix m. As you increase d, there is a higher prob. of points in general position being linearly separable.