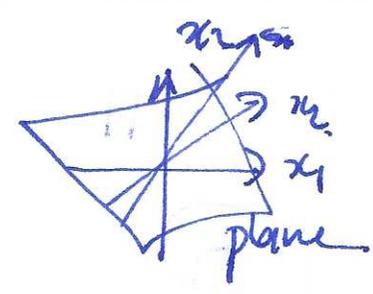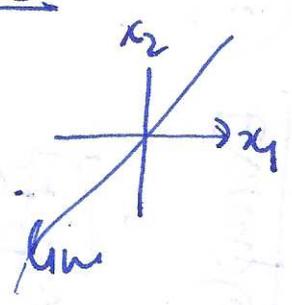# Support vector machines

hyperplane: $\mathbb{R}^{d-1} \subseteq \mathbb{R}^d$.

$$a_0 + a_1 x_1 + a_2 x_2 = 0$$

$$a_0 + a_1 x_1 + a_2 x_2 + a_3 x_3 = 0$$
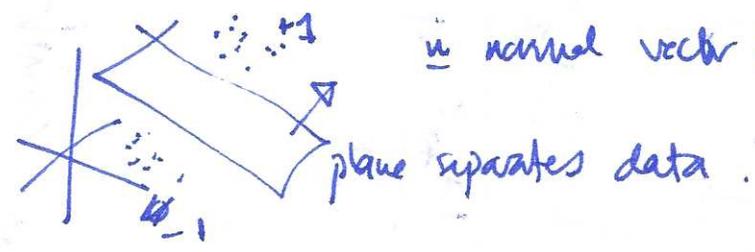
in $\mathbb{R}^d$:
$$a_0 + \sum_{i=1}^{d} a_i x_i = 0 .$$

$\underbrace{\qquad\qquad}$ function of $\underline{x} = (x_i)$

$\left. \begin{array}{l} f(\underline{x}) > 0 \\ f(\underline{x}) < 0 \end{array} \right\}$ two sides of hyperplane

## Classification.
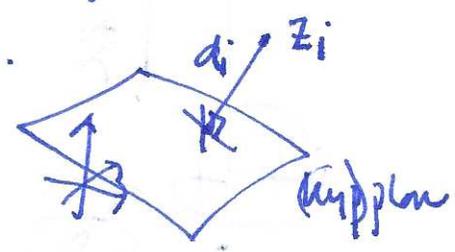
normal vector is $\underline{n} = (a_i)$.

plane separates data.

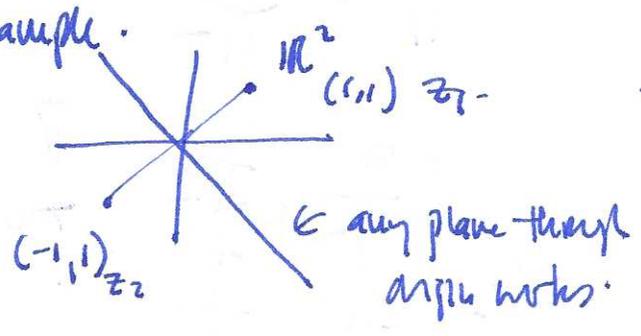get a bit more info $|f(x)| >> 0$ means point if far from the plane

## Maximal margin hyperplane / optimal separating hyperplane
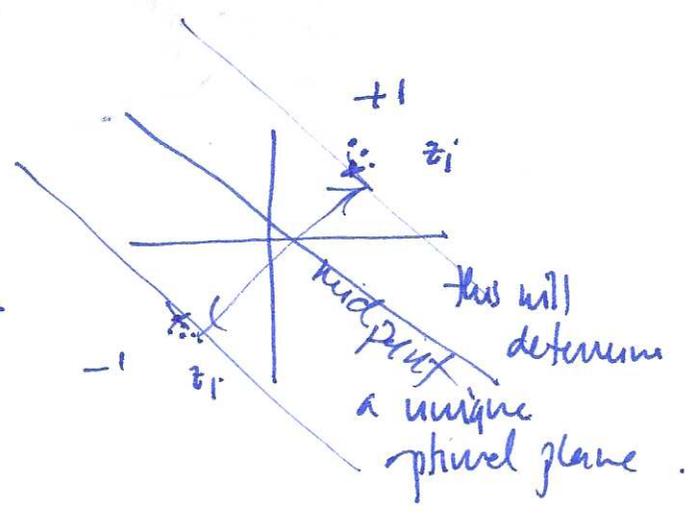
$Z = \{z_i\} \subseteq \mathbb{R}^d$.
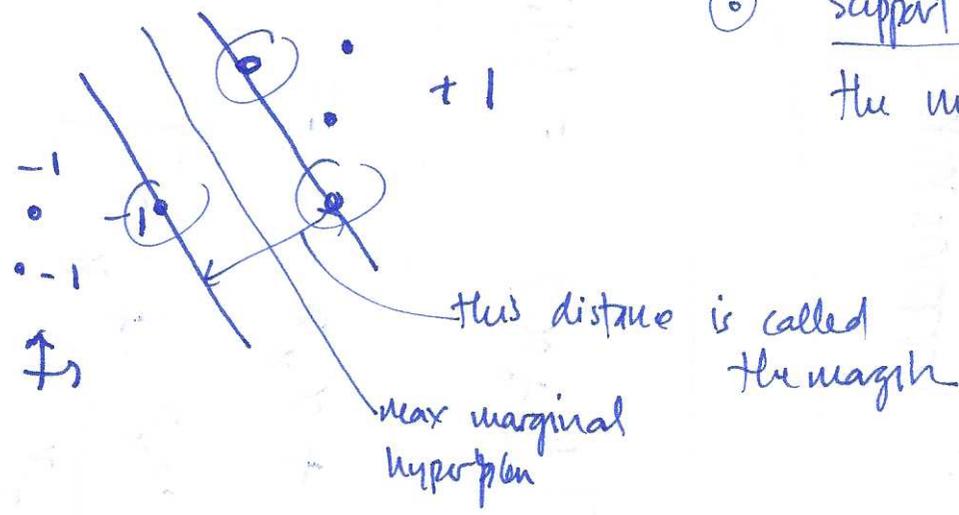
$d_i$ is called the margin.
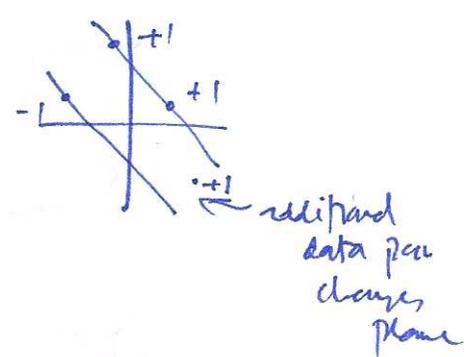
hyperplane

simplest example.

$\mathbb{R}^2$ $(1,1)$ $z_1$

$(-1,1)_{z_2}$

$\in$ any plane through origin works.

$+1$

$z_i$

this will determine a unique optimal plane.

(o) <u>support vectors</u>. these determine the maximal margin hyperplane!!

<u>unstable!</u>

this distance is called the margin

max marginal hyperplan

relocated data point changes plane

# Constructing the maximal margin classifier

setup $Z = \{z_i\} \subseteq \mathbb{R}^d$ labelled $\{\underset{y_i}{\pm 1}\}$. separated by a plane.

maximize $M$  subject to $\sum_{i=1}^{d} a_i^2 = 1$.
$\{a_i\}_{i=0}^{d}$   $\mathscr{L}(z)$.

$$y_i(a_0 + a_1 z_{i1} + a_2 z_{i2} + \cdots + a_d z_{id})$$

$$y_i(a_0 + a_1 z_{i1} + \cdots + a_d z_{id}) \geq M.$$

$$(\text{label}) \times f(z_i).$$

<u>Fact</u>: this optimization problem can be solved ∆.

<u>General case</u> $Z = \mathcal{A}(z_i)$ <u>not</u> separated by a hyperplane :

## Support Vector classifiers / soft margin classifier

$$\max_{\substack{a_i \, \epsilon_{j \leq n} \\ 1 \leq i \leq d}} M \quad (\text{set } \sum a_i^2 = 1)$$

subject to $\quad y_i(a_0 + a_1 z_{i1} \cdots a_d z_{id}) \geq M(1 - \epsilon_i)$

$$\epsilon_i \geq 0 \qquad \sum_{i=1}^{n} \epsilon_i < C$$

$\epsilon_i > 0$ means $z_i$ on wrong side of margin

$\epsilon_i > 1$          $z_i$ on wrong side of plane

+1   $\epsilon_i = 0$

+1   $\epsilon_i \geq 0$

+   $\epsilon_i \geq 1$   plane

$C = 0$ original case, ie. all points must be separated by hyperplane

$C$ tuning parameter ← chosen through cross validation.

in general:     $C$ small   plane classifier strongly fit to data

              loose   low bias high variance

         $C$ large    looser fit, more bias hopefully lower variance.

**key fact**    $z_i$ with $\epsilon_i = 0$ don't effect the choice of plane!

observe: the plane only depends on data points w/ $\epsilon_i > 0$

i.e. lying in the margin, or on the wrong side of the plane.

these are called the <u>support vectors</u>.

## <u>Support vector machines</u>.

aim: deal with non-linear boundaries.

$$\begin{array}{cc} -1 & -1 \\ -1 \quad +1 & -1 \\ -1 & -1 \end{array}$$

**fact**: finding the plane depends only on the <u>inner products</u>. of the data points

$\underline{a}\,\underline{b}. = (a_1, a_2 \dots, a_d) \cdot (b_1, \dots, b_d) = a_1 b_1 + a_2 b_2 + \dots + a_d b_d.$

• the linear support vector classifier can be represented as

$$f(x) = a_0 + \sum_{i=1}^{n} a_i \langle x, x_i \rangle \qquad \text{n parameters not 1 !}$$

to find the $a_i$ need the $\binom{n}{2}$ inner products $\langle x_i, x_i' \rangle$ for $(x_i) \in X$.

$$\frac{n(n-1)}{2} \sim n^2.$$

$a_i = 0$ unless $x_i$ is a support vector. Let $S$ be set of indices of support vectors.

$$f(x) = a_0 + \sum_{i \in S} \alpha_i \langle x, x_i \rangle$$

typically $|S| \ll \#\{n\}$

we can replace the inner product with ~~the~~ a _generalized inner product_ called a _kernel_ $K(x_i, x_i')$

standard inner product is $K(x_i, x_i) = x_i \cdot x_i' = \sum_j x_{ij} x_{ij}'$.

standard gives linear hyperplane classifier.

other choices : • $K(x_i, x_i') = \left( 1 + \sum_{j=1}^{d} x_{ij} x_{ij}' \right)^{k}$ &.

_polynomial kernel of degree $k$._

a _support vector machine_ is a classifier with a (possibly non-linear) kernel $K(x_i, x_i')$ and the classifying function will have the form $f(x) = a_0 + \sum_{i \in S} a_i K(x_i, x_i)$

• _radial kernel_ $K(x_i, x_i') = \exp\left( -\gamma \sum_{j=1}^{d} (x_{ij} - x_{ij}')^2 \right)$.

suppose there are more than two outcomes, not just $\{\pm 1\}$ ?

spse there are $\{1, 2, \ldots, K\}$ outcomes.

— can do $\binom{K}{2}$ pairwise comparisons and pick most frequent (one vs one).

— (one vs all) build $K$ SVMs identifying. $i$ against all other.

$\{1, \ldots K\} \setminus i \rightsquigarrow$ each of the $K$ classifiers gives a weight (+)

$f_i(x) \leftarrow$ choose biggest one.     margin

---

relation to logistic regression.

can rewrite criterion for finding the support vector classifier

$$f(x) = a_0 + a_1 x_1 + a_2 x_2 + \cdots + a_d x_d \quad \text{as}$$

$$\underset{a_i}{\text{minimize}} \left\{ \sum_{i=1}^{n} \max\{0, 1 - y_i f(x_i)\} + \lambda \sum_{j=1}^{d} a_i^2 \right\}$$

       ⊛           ↑

           non-negative tuning parameter

$\lambda$ large, $a_i$ are small, more margin violations, high bias low variance

$\lambda$ small    $a_i$ large       few margin violations   low bias, high variance

General setup : "Loss + Penalty"     ⌐ tuning parameter

$$\underset{a_i}{\text{minimize}} \left\{ L(X, y, a_i) + \lambda P(a_i) \right\}$$

        ↗                  ↑

    Loss function / error metric    penalty function on parameters.

measures how closely model fits data

Lasso :    $L(X, y, a) = \sum_{i=1}^{n} \left( y_i - a_0 - \sum_{j=1}^{d} x_{ij} a_j \right)^2$

                       ⊛⊛

$$P(a) = \sum_{i=1}^{n} |a_i| .$$

use ⊛ zero for non-support element

⊛⊛ very small for far away elements.