# SVM

Training data: $(\vec{z}_i, y_i)$ ; $i = 1 \ldots n$



+1 labelled

$w * x - b = 1$

$\frac{2}{\|w\|}$

$w * x - b = -1$

$\vec{w}$

O

$-1$ labelled.

Linearly separable data

Anything above $\vec{w}^T \vec{x} - b = 1$ is a "x" datapoint

Maximize $\frac{2}{\|\vec{w}\|}$ $\Rightarrow$ Minimize $\|w\|_2$

$\Rightarrow$ If $y_i = 1$, $\vec{w}^T \vec{z}_i - b \geq 1$ $\Big\}$ Points must be correctly
$\phantom{\Rightarrow If}$ $y_i = -1$, $\vec{w}^T \vec{z}_i - b \leq -1$. $\Big\}$ classified

$\Rightarrow$ $y_i (\vec{w}^T \vec{z}_i - b) \geq 1$ $\forall$ $1 \leq i \leq n$ $\qquad$ <u>Constraint</u>

Optimization problem:
$$\min_{w,b} \|\vec{w}\|_2^2$$

Quadratic optimization, solve w/ Lagrange multipliers
$\hookrightarrow$ Quadratic $\Rightarrow$ $\boxed{\text{single global minima}}$

If we remove support vectors, hyperplane changes!

sub. ~~Yilz~~ $y_i (\vec{w}^T \vec{z}_i - b) \geq 1$ $\forall i \in \{1 \ldots n\}$.

$\uparrow$ hyperplane is determined <u>completely</u> by nearest $z_i$ on either side!
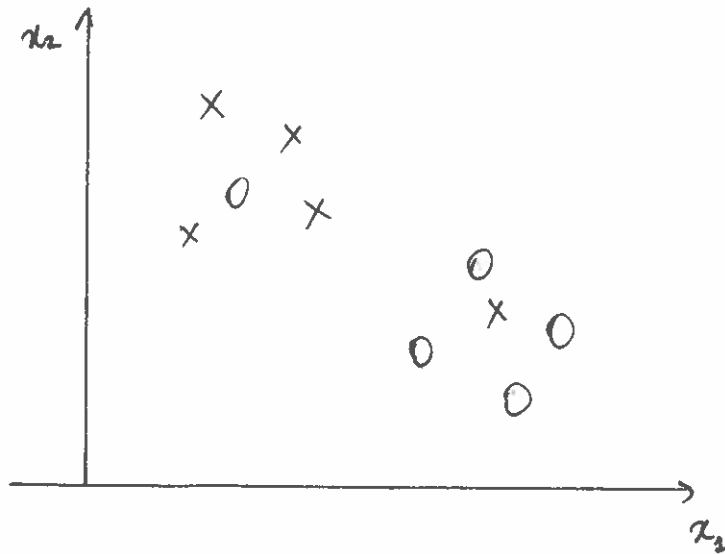
$\Rightarrow$ Most difficult pts to classify

We use the optimization of maximizing margin to reduce number of nonzero weights ~~that~~ $\Rightarrow$ only look at weights that matter $\Rightarrow$ ie. correspond to the support vectors!

Key diff b/w SVM & NNs/linear reg $\Rightarrow$ we all data!

Ultimately, classifier: $\text{sign}(\vec{w}^T\vec{z} - b)$ for some new $z$, using learned/ optimized $\vec{w}, b$.

"Fuzzy" data



Some points' labels are "misclassified"

Modify minimization:  (using hinge loss)

$$\max\left(0,\ 1 - y_i(\vec{w}^T\vec{z}_i - b)\right)$$

returns 0 if point is correctly classified/labelled $(0,1)$, proportional to distance from margin if incorrectly classified.

$\Rightarrow$ minimize $\lambda\|\vec{w}\|_2^2 + \left[\frac{1}{n}\sum_{i=1}^{n}\max\left(0,\ 1 - y_i(\vec{w}^T\vec{z}_i - b)\right)\right]$
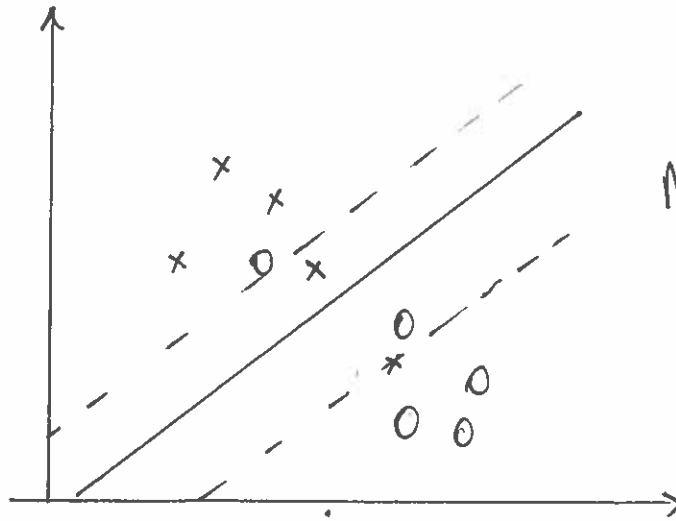
Rewrite: $\min\limits_{w,b,\zeta} \|\vec{w}\|_2^2 + C\sum_{i=1}^{n}\zeta_i$    Riemann zeta function

sub. $y_i(\vec{w}^T\vec{z}_i - b) \geq 1 - \zeta_i,\ \ \zeta_i \geq 0 \ \ \forall\, i \in \{1\ \dots\ n\}.$
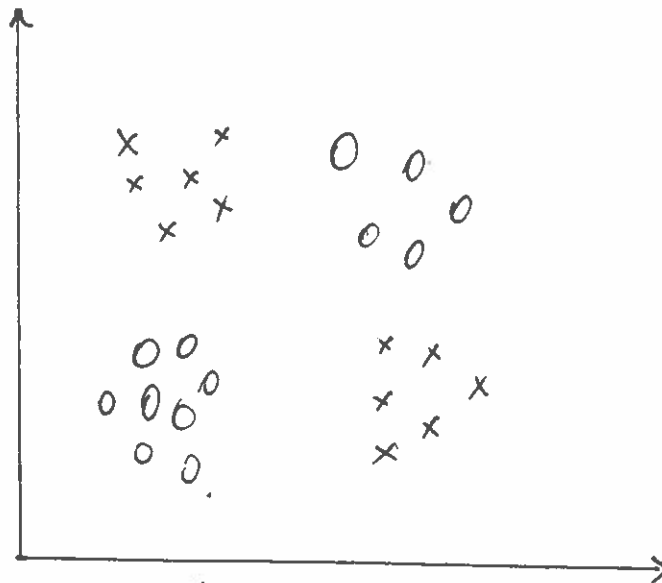
⇒ Our learned classifier now looks like:



Number of support vectors = 4

Very similar to linear discriminant analysis (LDA)

## Non-linearly Separable Data:



XOR dataset

Gain linear separation by mapping data into a high-dimensional space

We actually solve SVM optimization via the dual formulation (removes dependence on $\vec{w}, b$)

Dual problem –

Interchange variables & constraints :-

min $c \circ x$

sub. $A \circ x = b$

$x \geq 0$.

$\longrightarrow$

max $b \circ y$

sub. $A_i^T y \leq C$.

$y$ Free. (Unconst).

Here:

min $\|w\|_2^2 \longrightarrow$ | min. $f(\vec{w})$

sub. $y_i \cdot (\vec{w}^T \vec{x_i} - b) \ldots$ | sub. $g_j(\vec{w}) \leq 0$.

$h_k(\vec{w}) = 0$

General form ( Lin sep. / fuzzy).

| Dual formulation

$$\mathcal{L} = \frac{1}{2} f(\vec{w}) + \sum_j a_j g_j(\vec{w}) + \sum_k \beta_k h_k(\vec{w}) \nearrow 0 \because \text{no equalty const.}$$

(Method of Lagrange multipliers)

$\Rightarrow$ For a min /max:

$$\frac{\partial \mathcal{L}}{\partial w_i} = 0 \quad ; \quad \frac{\partial \mathcal{L}}{\partial \beta_k} = 0$$

From strong duality & Kahn-Tucker thms.,

max. $\mathcal{L}(\alpha, \beta)$.

sub. $a_j \geq 0 \;\; \forall j$. $\equiv$ original primal form.

Generally true $\forall$ convex frs. $\underline{f(w)}$.

Further, if $\hat{w}$ is optimal soln of primal & $\hat{a}$ & $\hat{\beta}$ are optimal

solns of dual:

$$f(\hat{w}) = \mathcal{L}(\hat{a}, \hat{b})$$
$$\hat{a}_j g_j(\hat{w}) = 0 \quad \forall j$$

} Karush-~~Kuter~~ Kuhn-Tucker (KKT) complementarity condition.

If we sub. our original constraints for the SVM problem into the Lagrangian,

$$\mathcal{L} = \frac{1}{2}\,\vec{w}\cdot\vec{w} + \sum_i \alpha_i \left[1 - y_i(\vec{w}\cdot\vec{x}_i + b)\right]$$

For max/min:

$$0 \stackrel{set}{=} \frac{\partial \mathcal{L}}{\partial \vec{w}} = \vec{w} - \sum_i \alpha_i y_i \vec{x}_i \qquad \text{and} \qquad 0 = \frac{\partial \mathcal{L}}{\partial b} = \sum_i \alpha_i y_i$$

$$\Rightarrow \hat{w} = \sum_i \hat{\alpha}_i y_i \vec{x}_i$$

$$\Rightarrow \mathcal{L}(\alpha) = \sum_j \alpha_j - \frac{1}{2}\sum_{j,k} \alpha_j y_j (\vec{x}_j \cdot \vec{x}_k) y_k \alpha_k$$
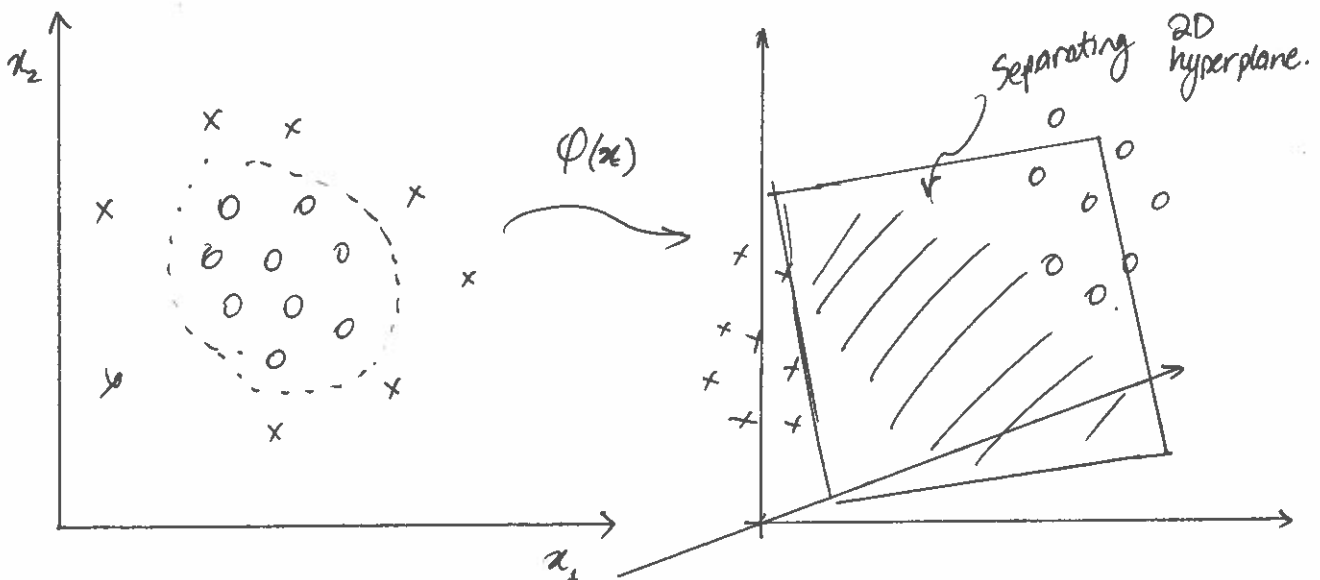
The only place where $\vec{x}$ shows up is here!
This is the only thing that scales w/ number of features. ⇒
all other parts scale w/ num. pts. ⇒ Primal to dual → shift scaling from
number (primal) of features to number of pts (dual!)
⇒ Favors data w/ ↓ numbers of data pts., but huge features

Now, lets talk about non linearity. ⇒ "Kernel Trick".



$\varphi(x)$

Separating 2D hyperplane.

**Primal form:** $\quad\min. \frac{1}{2}\vec{W}\cdot\vec{W} + \lambda\sum_i \Xi_i$

$$\text{sub.}\quad y_i(\vec{W}\cdot\varphi(\vec{x_i}) + B) \geq 1 - \Xi \qquad i=1\dots m.$$

But, embedding dim may be $\mathbb{R}^{Billions}$ !

$$\implies \text{Intractable!}$$

**Dual problem:** $\quad$ ~~$\min. \frac{1}{2}\alpha\cdot\text{diag}(y)$~~

$$\min. -\sum_j \alpha_j + \frac{1}{2}\sum_{j,k}\alpha_j y_j K_{jk} y_k \alpha_k$$

$$K_{jk} = K(x_j, x_k) = \varphi(x_i) \overset{\text{dot prod.}}{\cdot} \varphi(x_k)$$

But, we don't need to explicitly know mapping $\varphi(x)$ !

All we need is some way to compute $K_{jk}$ that could have come from some $\varphi(x)$ ! $\implies$ $m\times m$ matrix w/ mathematical properties of inner product!

- $K_{ij} = K(x_i, x_j)$ must be symm. in $i, j$, and non neg. eigenval.
- $K$ can be composed by addition, mult., scaling by const.

**Popular Kernels:**

$$\text{linear:}\quad K(x_i, x_j) = x_i \cdot x_j.$$

$$\text{power:}\quad K(x_i, x_j) = (x_i \cdot x_j)^d \qquad 2 \leq d \leq 20 \text{ (usually)}.$$

$$\text{polynom.:}\quad K(x_i, x_j) = (a\, x_i \cdot x_j + b)^d$$

$$\text{sigmoid:}\quad K(x_i, x_j) = \tanh(a x_i \cdot x_j + b).$$

$$\text{Gaussian rbf:}\quad K(x_i, x_j) = e^{-\frac{1}{2}\frac{|x_i - x_j|^2}{\sigma^2}}$$

Tips for using kernels:

- Gaussian rbf is very popular, ∴ only 1 hyperparam.
   Guess good initial $\sigma$ via avg. distance b/w points in <u>feature space.</u>

- Polynomial kernels:

   Start by choosing $a, b$ S.T. $a x_i \cdot x_j + b$ b/w $-1, 1$ $\forall i, j$.

   $d \Rightarrow$ roughly interpret as number of features to be mixed ~~to~~
      during partitioning.

               $\Rightarrow d=1 \Rightarrow$ partition space by 1 feature at a time
               $2 \Rightarrow$   2 features ...


   Diff b/w power & polynomial: power considers only $d$ features at
      once, polynomial considers all combos of $d$ or fewer features.