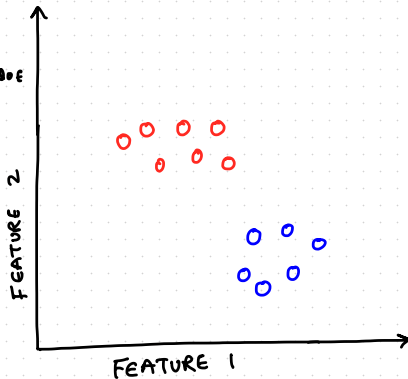# Support Vector Machines

Nipun Batra

April 23, 2023

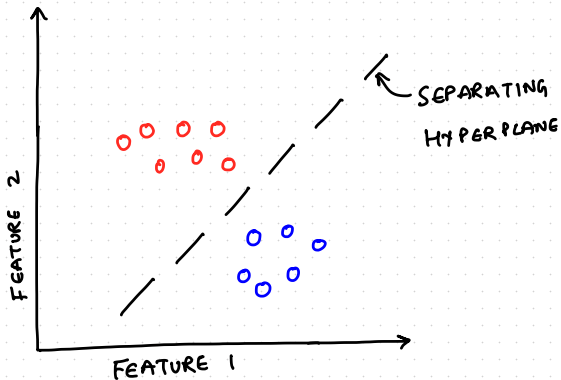IIT Gandhinagar

# SUPPORT VECTOR MACHINES
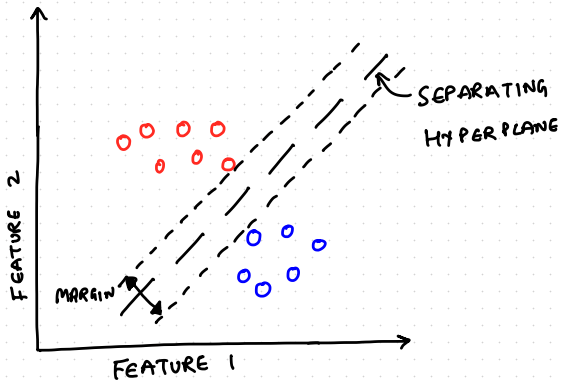
## POPULAR BINARY

CLASSIFICATION TECHNIQUE



FEATURE 2

FEATURE 1

IDEA: DRAW A SEPARATING HYPERPLANE

FEATURE 2

FEATURE 1

SEPARATING HYPERPLANE

MARGIN

IDEA: MAXIMIZE THE MARGIN
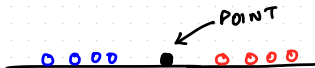
SEPARATING HYPERPLANE

SUPPORT VECTORS

FEATURE 2

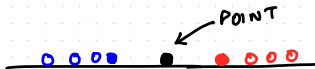MARGIN

FEATURE 1

SUPPORT VECTORS: POINTS ON BOUNDARY | MARGIN

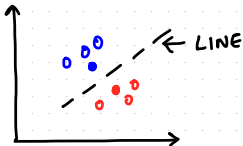# HYPERPLANE VIS # DIMENSIONS

1D

POINT

# HYPERPLANE VS # DIMENSIONS

1D

POINT

2D

LINE

# Hyperplane vs # Dimensions

## 1D

POINT

## 2D

← LINE

## 3D (and more)

← HYPERPLANE

WHICH   HYPERPLANE?

FEATURE 2

FEATURE 1

WHICH HYPERPLANE?

# EQUATION OF HYPERPLANE

HOW TO DEFINE?

ORIGIN

# EQUATION OF HYPERPLANE

$P$ : Any point on plane

$P_0$ : One point on plane



·$P$   · $P_0$

ORIGIN

# EQUATION OF HYPERPLANE



$\vec{w}$ : $\perp$ vector to plane at $P_0$

# EQUATION OF HYPERPLANE



$\vec{w}$

P

$P_0$

$\vec{x}$

$\vec{x}_0$

ORIGIN

P and $P_0$ lie on plane

# EQUATION OF HYPERPLANE



$\vec{PP_0} = \vec{x} - \vec{x_0}$ lies on plane

# EQUATION OF HYPERPLANE



$\vec{P P_0} = \vec{x} - \vec{x_0}$ lies on plane

$\Rightarrow \quad \vec{w} \perp (\vec{x} - \vec{x_0})$

or, $\quad \vec{w} \cdot (\vec{x} - \vec{x_0}) = 0$

or, $\quad \vec{w} \cdot \vec{x} - \vec{w} \cdot \vec{x_0} = 0$

or, $\quad \boxed{\vec{w} \cdot \vec{x} + b = 0}$

# DISTANCE B/W || HYPER PLANES



$\vec{w}.\vec{x} + b_2 = 0$

$\vec{w}.\vec{x} + b_1 = 0$

# DISTANCE B/W || HYPERPLANES



$\vec{w} \cdot \vec{x} + b_2 = 0$

$t\vec{w}$

$\vec{x_2}$

$\vec{x_1}$

$\vec{w} \cdot \vec{x} + b_1 = 0$

ORIGIN

## Distance between 2 parallel hyperplanes

Equation of two planes is:

$$\vec{w} \cdot \vec{x} + b_1 = 0$$
$$\vec{w} \cdot \vec{x} + b_2 = 0$$

## Distance between 2 parallel hyperplanes

Equation of two planes is:

$$\vec{w} \cdot \vec{x} + b_1 = 0$$

$$\vec{w} \cdot \vec{x} + b_2 = 0$$

For a point $\vec{x_1}$ on plane 1 and $\vec{x_2}$ on plane 2, we have:

## Distance between 2 parallel hyperplanes

Equation of two planes is:

$$\vec{w} \cdot \vec{x} + b_1 = 0$$
$$\vec{w} \cdot \vec{x} + b_2 = 0$$

For a point $\vec{x_1}$ on plane 1 and $\vec{x_2}$ on plane 2, we have:

$$\vec{x_2} = \vec{x_1} + t\vec{w}$$
$$D = |t\vec{w}| = |t|||\vec{w}||$$

### Distance between 2 parallel hyperplanes

Equation of two planes is:

$$\vec{w} \cdot \vec{x} + b_1 = 0$$
$$\vec{w} \cdot \vec{x} + b_2 = 0$$

For a point $\vec{x_1}$ on plane 1 and $\vec{x_2}$ on plane 2, we have:

$$\vec{x_2} = \vec{x_1} + t\vec{w}$$
$$D = |t\vec{w}| = |t| ||\vec{w}||$$

We can rewrite as follows:

### Distance between 2 parallel hyperplanes

Equation of two planes is:

$$\vec{w} \cdot \vec{x} + b_1 = 0$$
$$\vec{w} \cdot \vec{x} + b_2 = 0$$

For a point $\vec{x_1}$ on plane 1 and $\vec{x_2}$ on plane 2, we have:

$$\vec{x_2} = \vec{x_1} + t\vec{w}$$
$$D = |t\vec{w}| = |t|\,||\vec{w}||$$

We can rewrite as follows:

$$\vec{w} \cdot \vec{x_2} + b_2 = 0$$
$$\Rightarrow \vec{w} \cdot (\vec{x_1} + t\vec{w}) + b_2 = 0$$

## Distance between 2 parallel hyperplanes

Equation of two planes is:

$$\vec{w} \cdot \vec{x} + b_1 = 0$$
$$\vec{w} \cdot \vec{x} + b_2 = 0$$

For a point $\vec{x_1}$ on plane 1 and $\vec{x_2}$ on plane 2, we have:

$$\vec{x_2} = \vec{x_1} + t\vec{w}$$
$$D = |t\vec{w}| = |t| \|\vec{w}\|$$

We can rewrite as follows:

$$\vec{w} \cdot \vec{x_2} + b_2 = 0$$
$$\Rightarrow \vec{w} \cdot (\vec{x_1} + t\vec{w}) + b_2 = 0$$

$$\Rightarrow \vec{w} \cdot \vec{x_1} + t \|\vec{w}\|^2 + b_1 - b_1 + b_2 = 0 \Rightarrow t = \frac{b_1 - b_2}{\|\vec{w}\|^2} \Rightarrow D = t \|\vec{w}\| = \frac{b_1 - b_2}{\|\vec{w}\|}$$

FORMULATION

+1 CLASS

−1 CLASS

FEATURE 2

FEATURE 1

FORMULATION

$\vec{w} \cdot \vec{x} + b = +1$

$\vec{w} \cdot \vec{x} + b = -1$

FEATURE 2

FEATURE 1

FORMULATION

$\vec{w}.\vec{x}+b=+1$

$\vec{w}.\vec{x}+b=-1$

FEATURE 2

FEATURE 1

$$\text{MARGIN} = \frac{(b+1)-(b-1)}{||\vec{w}||}$$

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

# FORMULATION



**GOAL:** MAXIMIZE MARGIN

$\Rightarrow$ MAXIMIZE $\dfrac{2}{||\vec{w}||}$

$\Rightarrow$ MINIMIZE $||\vec{w}||$

S.T. Correctly label points

FORMULATION

GOAL: MAXIMIZE MARGIN

$\Rightarrow$ MAXIMIZE $\dfrac{2}{||\vec{w}||}$

$\Rightarrow$ MINIMIZE $||\vec{w}||$

S.T. Correctly label points

i.e. if $y_i = -1$

$\vec{w} \cdot \vec{x} + b \leq -1$

if $y_i = +1$

$\vec{w} \cdot \vec{x} + b \geq 1$

FORMULATION



GOAL: MAXIMIZE MARGIN

$\Rightarrow$ MAXIMIZE $\dfrac{2}{||\vec{w}||}$

$\Rightarrow$ MINIMIZE $||\vec{w}||$

S.T. Correctly label points

i.e. if $y_i = -1$

$\quad \vec{w} \cdot \vec{x} + b \leq -1$

if $y_i = +1$

$\quad \vec{w} \cdot \vec{x} + b \geq 1$

$y_i (\vec{w} \cdot \vec{x} + b) \geq 1$

## Primal Formulation

Objective

$$\text{Minimize } \frac{1}{2}||w||^2$$
$$\text{s.t. } y_i(w.x_i + b) \geq 1 \ \forall i$$

## Primal Formulation

> Objective
>
> $$\text{Minimize } \frac{1}{2}||w||^2$$
> $$\text{s.t. } y_i(w.x_i + b) \geq 1 \;\; \forall i$$

Q) What is $||w||$?

## Primal Formulation

> Objective
>
> $$\text{Minimize } \frac{1}{2}||w||^2$$
> $$\text{s.t. } y_i(w.x_i + b) \geq 1 \ \ \forall i$$

Q) What is $||w||$?

$$w = \begin{bmatrix} w_1 \\ w_2 \\ ... \\ w_n \end{bmatrix}$$

$$||w|| = \sqrt{w^T w}$$

$$= \sqrt{\begin{bmatrix} w_1, w_2, ... w_n \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ ... \\ w_n \end{bmatrix}}$$

# EXAMPLE (IN 1D)



SEPARATING POINT

## Simple Exercise

$$\begin{bmatrix} x & y \\ 1 & 1 \\ 2 & 1 \\ -1 & -1 \\ -2 & -1 \end{bmatrix}$$

Separating Hyperplane: $wx + b = 0$

## Simple Exercise

$$y_i(w_i x_i + b) \geq 1$$

$$\begin{bmatrix} x_1 & y \\ 1 & 1 \\ 2 & 1 \\ -1 & -1 \\ -2 & -1 \end{bmatrix}$$

$\Rightarrow y_i(w_i x_i + b) \geq 1$

$\Rightarrow 1(w_1 + b) \geq 1$

$\Rightarrow 1(2w_1 + b) \geq 1$

$\Rightarrow -1(-w_1 + b) \geq 1$

$\Rightarrow -1(-2w_1 + b) \geq 1$

$w + b = 1$

①

$w+b=1$
①

$b$

$3$

$2$

$1$

$\frac{1}{2}$   $1$   $2$   $3$   $\omega$

$2\omega+b=1$
②

$\omega+b=1$
①

$b$

$3$

$2$

$1$

$w$

$\frac{1}{2}$  $1$  $2$  $3$

$w-b=1$

③

② $2w+b=1$

① $w+b=1$

$b$

④ $2w+b=1$

③ $w-b=1$

$w$

$w+b=1$
①

② $2w+b=1$

## Simple Exercise

$$w_{min} = 1, b = 0$$
$$w.x + b = 0$$
$$x = 0$$

## Simple Exercise

Minimum values satisfying constraints $\Rightarrow w = 1$ and $b = 0$

$\therefore$ Max margin classifier $\Rightarrow x = 0$

## Primal Formulation is a Quadratic Program

Generally;

$\Rightarrow$ Minimize Quadratic(x)

$\Rightarrow$ such that, Linear(x)

Question

$$x = (x_1, x_2)$$

minimize $\dfrac{1}{2}||x||^2$

$: x_1 + x_2 - 1 \geq 0$

MINIMIZE  QUADRATIC

s.t.  LINEAR



$x_2$

SOLUTION

$x_1$

$x_1 + x_2 - 1 \geq 0$

## Converting to Dual Problem

Primal $\Rightarrow$ Dual Conversion using Lagrangian multipliers

$$\text{Minimize } \frac{1}{2}||\bar{w}||^2$$
$$\text{s.t. } y_i(\bar{w}.x_i + b) \geq 1$$
$$\forall i$$

$$L(\bar{w}, b, \alpha_1, \alpha_2, ...\alpha_n) = \frac{1}{2}\sum_{i=1}^{d} w_i^2 - \sum_{i=1}^{N} \alpha_i(y_i(\bar{w}.\bar{x}_i + b) - 1) \ \forall \ \alpha_i \geq 0$$

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

## Converting to Dual Problem

$$\frac{\partial L}{\partial w} = 0 \Rightarrow \bar{w} - \sum_{i=1}^{n} \alpha_i y_i \bar{x}_i = 0$$

$$\bar{w} = \sum_{i=1}^{N} \alpha_i y_i \bar{x}_i$$

$$L(\bar{w}, b, \alpha_1, \alpha_2, ...\alpha_n) = \frac{1}{2} \sum_{i=1}^{d} w_i^2 - \sum_{i=1}^{N} \alpha_i (y_i(\bar{w}.\bar{x}_i + b) - 1$$

$$= \frac{1}{2} ||\bar{w}||^2 - \sum_{i=1}^{N} \alpha_i y_i \vec{w}.\bar{x}_i - \sum_{i=1}^{N} \alpha_i y_i b + \sum_{i=1}^{N} \alpha_i$$

$$= \sum_{=1}^{N} \alpha_i + \frac{\left(\sum_i \alpha_i y_i \bar{x}_i\right)\left(\sum_j \alpha_j y_j \bar{x}_j\right)}{2} - \sum_i \alpha_i y_i \left(\sum_j \alpha_j y_j \bar{x}_j\right) \bar{x}_i$$

## Converting to Dual Problem

$$L(\alpha) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j \bar{x}_i \cdot \bar{x}_j$$

Minimize $\|\bar{w}\|^2 \Rightarrow$    Maximize $L(\alpha)$

s.t                 s.t

$y_i (\bar{w}, x_i + b) \geqslant 1$      $\sum_{i=1}^{N} \alpha_i y_i = 0 \ \forall \ \alpha_i \geq 0$

## Question

**Question**:

$\alpha_i \left( y_i \left( \bar{w}, \bar{x}_i + b \right) - 1 \right) = 0 \quad \forall i$ as per KKT slackness

What is $\alpha_i$ for support vector points?

**Answer:** For support vectors,

$$\bar{w}.\bar{x}_i + b = -1 \text{ (+ve class)}$$
$$\bar{w}.\bar{x}_i + b = +1 \text{ (+ve class)}$$

$y_i \left( \bar{w} \cdot \bar{x}_i + b \right) - 1 \right) = 0 \quad$ for $i = \{\text{support vector points}\}$

$\therefore \alpha_i$ where i $\in \{\text{support vector points}\} \neq 0$

For all non-support vector points $\alpha_i = 0$

EXAMPLE (IN 1D)



SEPARATING POINT

## Revisiting the Simple Example

$$\begin{bmatrix} x_1 & y \\ 1 & 1 \\ 2 & 1 \\ -1 & -1 \\ -2 & -1 \end{bmatrix}$$

$$L(\alpha) = \sum_{i=1}^{4} \alpha_i - \frac{1}{2} \sum_{i=1}^{4} \sum_{j=1}^{4} \alpha_i \alpha_j y_i y_j \bar{x}_i \bar{x}_j \qquad \alpha_i \geq 0$$

$$\sum \alpha_i y_i = 0 \qquad \alpha_i (y_i (\bar{w}.\bar{x}_i + b - 1) = 0$$

### Revisiting the Simple Example

$$
\begin{aligned}
L(\alpha_1, \alpha_2, \alpha_3, \alpha_4) = &\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4 \\
&- \frac{1}{2} \{ \alpha_1 \alpha_1 \times (1 * 1) \times (1 * 1) \\
&\quad + \\
&\alpha_1 \alpha_2 \times (1 * 1) \times (1 * 2) \\
&\quad + \\
&\alpha_1 \alpha_3 \times (1 * -1) \times (1 * 1) \\
&\quad ... \\
&\alpha_4 \alpha_4 \times (-1 * -1) \times (-2 * -2) \}
\end{aligned}
$$

How to Solve? $\Rightarrow$ Use the QP Solver!!

## Revisiting the Simple Example

For the trivial example,

We know that only $x = \pm 1$ will take part in the constraint actively.

Thus, $\alpha_2, \alpha_4 = 0$

By symmetry, $\alpha_1 = \alpha_3 = \alpha$ (say)

& $\sum y_i \alpha_i = 0$

$$L(\alpha_1, \alpha_2, \alpha_3, \alpha_4) = 2\alpha$$
$$- \frac{1}{2} \left\{ \alpha^2(1)(-1)(1)(-1) \right.$$
$$+ \alpha^2(-1)(1)(-1)(1)$$
$$\left. + \alpha^2(1)(1)(1)(1) + \alpha^2(-1)(-1)(-1)(-1) \right\}$$

$$\underset{\alpha}{Maximize} \quad 2\alpha - \frac{1}{2}(4\alpha^2)$$

14

## Revisiting the Simple Example

$$\frac{\partial}{\partial \alpha} \left(2\alpha - 2\alpha^2\right) = 0 \Rightarrow 2 - 4\alpha = 0$$

$$\Rightarrow \alpha = 1/2$$

$$\therefore \alpha_1 = 1/2 \ \ \alpha_2 = 0; \ \ \alpha_3 = 1/2 \ \ \alpha_4 = 0$$

$$\vec{w} = \sum_{i=1}^{N} \alpha_i y_i \bar{x}_i = 1/2 \times 1 \times 1 + 0 \times 1 \times 2$$

$$+ 1/2 \times -1 \times -1 + 0 \times -1 \times -2$$

$$= 1/2 + 1/2 = 1$$

**Finding b:**

For the support vectors we have,

$y_i(\vec{w} \cdot \overrightarrow{x_i} + b) - 1 = 0$

or, $y_i (\bar{w} \cdot \bar{x}_1 + b) = 1$

or, $y_i^2 (\bar{w} \cdot \bar{x}_i + b) = y_i$

or, $\bar{w}, \bar{x}_i + b = y_i \ (\because y_i^2 = 1)$

or, $b = y_i - w \cdot x_i$

In practice, $b = \frac{1}{N_{SV}} \sum_{i=1}^{N_{SV}} (y_i - \bar{w}\bar{x}_i)$

## Obtaining the Solution

$$b = \frac{1}{2}\{(1 - (1)(1)) + (-1 - (1)(-1))$$
$$= \frac{1}{2}\{0 + 0\} = 0$$
$$= 0$$
$$\therefore w = 1 \ \& \ b = 0$$

**Making Predictions**
$$\hat{y}(x_i) = \text{SIGN}(w \cdot x_i + b)$$
For $x_{test} = 3$; $\hat{y}(3) = \text{SIGN}(1 \times 3 + 0) = +\text{ve class}$
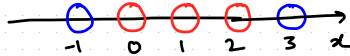
## Making Predictions

Alternatively,
$$\hat{y}\left(x_{TEST}\right) = \mathsf{SIGN}\left(\bar{w} \cdot \bar{x}_{TEST} + b\right)$$
$$= \mathsf{SIGN}\left(\sum_{i=1}^{N_S} \alpha_j y_j x_j \cdot x_{test} + b\right)$$

In our example,

$\alpha_1 = 1/2; \alpha_2 = 0; \quad \alpha_3 = 1/2; \alpha_4 = 0$

$\hat{y}(3) = \mathsf{SIGN}\left(\dfrac{1}{2} \times 1 \times (1 \times 3) + 0 + \dfrac{1}{2} \times (-1) \times (-1 \times 3) + 0\right)$

$= \mathsf{SIGN}\left(\dfrac{6}{2}\right) = \mathsf{SIGN}(3) = +1$

ORIGINAL DATA
IN R

# Non-Linearly Separable Data

Data not separable in $\mathbb{R}$

Data not separable in $\mathbb{R}$
Can we still use SVM?

## Non-Linearly Separable Data

Data not separable in $\mathbb{R}$
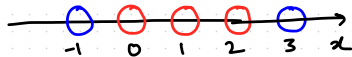Can we still use SVM?
Yes!

## Non-Linearly Separable Data
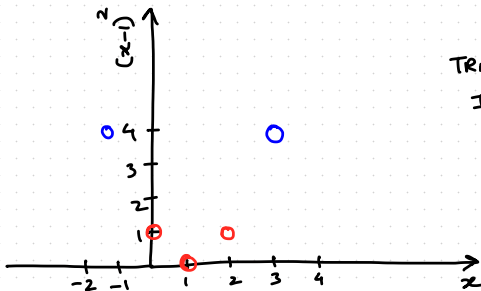
Data not separable in $\mathbb{R}$
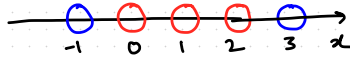
Can we still use SVM?
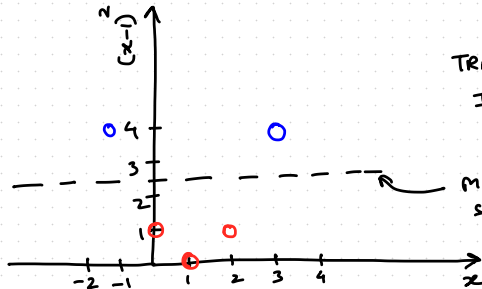
Yes!

How? Project data to a higher dimensional space.

ORIGINAL DATA
IN R

TRANSFORMED DATA
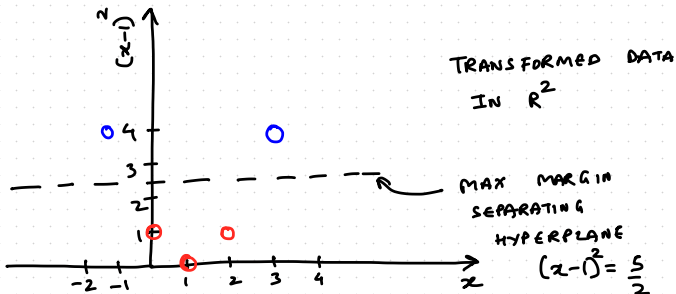IN R²

ORIGINAL DATA
IN $\mathbb{R}$

TRANSFORMED DATA
IN $\mathbb{R}^2$

MAX MARGIN
SEPARATING
HYPERPLANE

$(x-1)^2 = \dfrac{5}{2}$

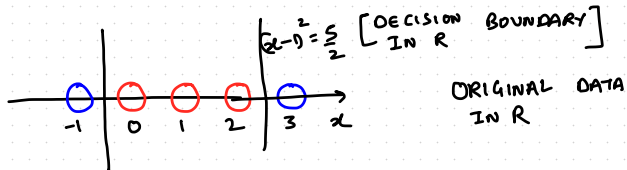$(x-1)^2 = \frac{5}{2}$ [DECISION BOUNDARY IN $\mathbb{R}$]
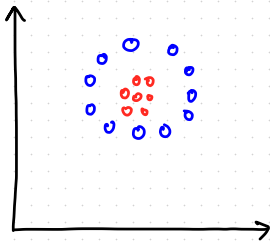
ORIGINAL DATA IN $\mathbb{R}$

TRANSFORMED DATA IN $\mathbb{R}^2$

MAX MARGIN SEPARATING HYPERPLANE

$(x-1)^2 = \frac{5}{2}$

$R^2$ SPACE

$R^2$ SPACE

$\Rightarrow$

$R^3$ SPACE

DECISION BOUNDARY

$R^2$ SPACE

DECISION BOUNDARY IN ORIGINAL SPACE

$R^3$ SPACE

DECISION BOUNDARY

ORIGINAL DATA IN R

TRANSFORMED DATA IN $R^2$

$$\phi(x) = \begin{bmatrix} \sqrt{2}\,x \\ x^2 \end{bmatrix}$$

ORIGINAL DATA
IN R

TRANSFORMED DATA
IN $R^2$

$$\phi(x) = \begin{bmatrix} \sqrt{2}\,x \\ x^2 \end{bmatrix}$$

ORIGINAL DATA
IN R

TRANSFORMED DATA
IN $R^2$

$$\phi(x) = \begin{bmatrix} \sqrt{2}\,x \\ x^2 \end{bmatrix}$$

$x_2$

$(0,1)$

$(-1,0)$      $(1,0)$    $x_1$

$(0,-1)$

$x_2$
$(= x_2^2)$

$x_1 (= x_1^2)$

$x_3 (= \sqrt{2}\, x_1 x_2)$

$x_2$

$(0,1)$

$(-1,0)$          $(1,0)$          $x_1$

$(0,-1)$

$x_2$
$(= x_2^2)$

SEPARATING        HYPER PLANE

$x_1 (= x_1^2)$

$x_3 (= \sqrt{2}\, x_1 x_2)$

## Linear SVMs in higher dimensions

Linear SVM:

Maximize

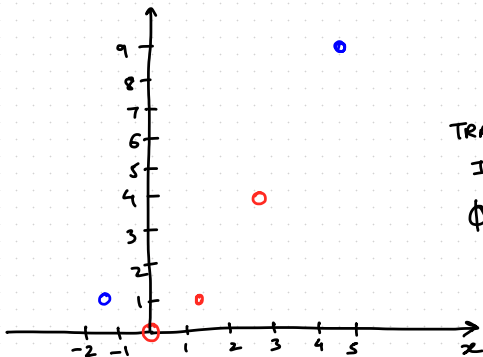$$L(\alpha) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j \overline{x_i}.\overline{x_j}$$

such that constriants are satisfied.

$$\downarrow$$
Transformation $(\phi)$
$$\downarrow$$

$$L(\alpha) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j \phi(\overline{x_i}).\phi(\overline{x_j})$$

**Linear SVMs in higher dimensions: Steps**

1. Compute $\phi(x)$ for each point

$$\phi : \mathbb{R}^d \to \mathbb{R}^D$$

2. Compute dot products over $\mathbb{R}^D$ space

Q. If $D >> d$

Both steps are expensive!

# Kernel Trick

## Kernel Trick

- Can we compute $K(\bar{x}_i, \bar{x}_j)$ , such that

- Can we compute $K(\bar{x}_i, \bar{x}_j)$ , such that
- $K(\bar{x}_i, \bar{x}_j) = \phi(\bar{x}_i).\phi(\bar{x}_j)$ , where

## Kernel Trick

- Can we compute $K(\bar{x}_i, \bar{x}_j)$ , such that
- $K(\bar{x}_i, \bar{x}_j) = \phi(\bar{x}_i).\phi(\bar{x}_j)$ , where
- $K(\bar{x}_i, \bar{x}_j)$ is some function of dot product in original dimension

## Kernel Trick

- Can we compute $K(\bar{x}_i, \bar{x}_j)$ , such that
- $K(\bar{x}_i, \bar{x}_j) = \phi(\bar{x}_i).\phi(\bar{x}_j)$ , where
- $K(\bar{x}_i, \bar{x}_j)$ is some function of dot product in original dimension
- $\phi(\bar{x}_i).\phi(\bar{x}_j)$ is dot product in high dimensions (after transformation)

KERNEL

TRICK

ORIGINAL DATA
IN R

TRANSFORMED DATA
IN $R^2$

$$\phi(x) = \begin{bmatrix} \sqrt{2}\,x \\ x^2 \end{bmatrix}$$

# KERNEL TRICK

$$\phi(x) = \begin{bmatrix} \sqrt{2}\,x \\ x^2 \end{bmatrix}$$

$$K(x_i, x_j) = \quad ?$$

# KERNEL TRICK

$$\phi(x) = \begin{bmatrix} \sqrt{2}\,x \\ x^2 \end{bmatrix}$$

$$K(x_i, x_j) = (1 + x_i \cdot x_j)^2 - 1$$

# KERNEL TRICK

$$\phi(x) = \begin{bmatrix} \sqrt{2}\,x \\ x^2 \end{bmatrix}$$

$$K(x_i, x_j) = (1 + x_i \cdot x_j)^2 - 1$$

$$(1 + x_i \cdot x_j)^2 - 1 = 1 + 2\,x_i \cdot x_j + x_i^2 x_j^2 - 1$$

$$= 2\,x_i \cdot x_j + x_i^2 x_j^2$$

$$= \left( \sqrt{2}\,x_i \cdot \sqrt{2}\,x_j + x_i^2 \cdot x_j^2 \right)$$

$$= \langle \sqrt{2}\,x_i, x_i^2 \rangle \cdot \langle \sqrt{2}\,x_j, x_j^2 \rangle$$

$$= \phi(x_i) \cdot \phi(x_j)$$

ORIGINAL   DATASET

| $x$ | $y$ |
|-----|-----|
| -1  | -1  |
| 0   | 1   |
| 1   | 1   |
| 2   | 1   |
| 3   | -1  |

ORIGINAL DATASET

| $x$ | $y$ |
|---|---|
| -1 | -1 |
| 0 | 1 |
| 1 | 1 |
| 2 | 1 |
| 3 | -1 |

TRANSFORMED DATASET

| $x$ | $\sqrt{2}x$ | $x^2$ | $y$ |
|---|---|---|---|
| -1 | $-\sqrt{2}$ | 1 | -1 |
| 0 | 0 | 0 | 1 |
| 1 | $\sqrt{2}$ | 1 | 1 |
| 2 | $2\sqrt{2}$ | 4 | 1 |
| 3 | $3\sqrt{2}$ | 9 | -1 |

ORIGINAL DATASET

| $x$ | $y$ |
|---|---|
| -1 | -1 |
| 0 | 1 |
| 1 | 1 |
| 2 | 1 |
| 3 | -1 |

TRANSFORMED DATASET

| $x$ | $\sqrt{2}x$ | $x^2$ | $y$ |
|---|---|---|---|
| -1 | $-\sqrt{2}$ | 1 | -1 |
| 0 | 0 | 0 | 1 |
| 1 | $\sqrt{2}$ | 1 | 1 |
| 2 | $2\sqrt{2}$ | 4 | 1 |
| 3 | $3\sqrt{2}$ | 9 | -1 |

Calculation w/o Kernel Trick

$$\phi(x_1) = \langle \sqrt{2}\,x, x^2 \rangle : 2$$

$$\phi(x_2) = \langle \sqrt{2}\,x, x^2 \rangle : 2$$

$$\phi(x_1) \cdot \phi(x_2) = 2 \text{ MULTIPLICATION} + 1 \text{ Addition}$$

ORIGINAL DATASET

| $x$ | $y$ |
|-----|-----|
| -1 | -1 |
| 0 | 1 |
| 1 | 1 |
| 2 | 1 |
| 3 | -1 |

TRANSFORMED DATASET

| $x$ | $\sqrt{2}x$ | $x^2$ | $y$ |
|-----|-----|-----|-----|
| -1 | $-\sqrt{2}$ | 1 | -1 |
| 0 | 0 | 0 | 1 |
| 1 | $\sqrt{2}$ | 1 | 1 |
| 2 | $2\sqrt{2}$ | 4 | 1 |
| 3 | $3\sqrt{2}$ | 9 | -1 |

Calculation with Kernel Trick

$$K(x_1, x_2) = \left(1 + x_1 \cdot x_2\right)^2 - 1$$

$x \cdot x_2 \longrightarrow 1$

$1 + x_1 \cdot x_2 \longrightarrow 1$

$\left(1 + x_1 \cdot x_2\right)^2 \longrightarrow 1$

$\left(1 + x_1 \cdot x_2\right)^2 - 1 \longrightarrow 1$

$\Big\} \; 4$

## Kernel Trick

Q) Why did we use dual form?

## Kernel Trick

Q) Why did we use dual form?
Kernels again!!

## Kernel Trick

Q) Why did we use dual form?

Kernels again!!

Primal form doesn't allow for the kernel trick $K(\bar{x}_1, \bar{x}_2)$ in dual and compute $\phi(x)$ and then dot product in $D$ dimensions

## Some Kernels

1. Linear: $K(\bar{x}_1, \bar{x}_2) = \bar{x}_1 \bar{x}_2$
2. Polynomial: $K(\bar{x}_1, \bar{x}_2) = (p + \bar{x}_1 \bar{x}_2)^q$
3. Gaussian: $K(\bar{x}_1, \bar{x}_2) = e^{-\gamma \|\bar{x}_1 - \bar{x}_2\|^2}$ where $\gamma = \frac{1}{2\sigma^2}$ - Also called Radial Basis Function (RBF)

## Kernels

Q) For $\bar{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ what space does kernel $K(\bar{x}, \bar{x'}) = (1 + \bar{x}\bar{x'})^3$

belong to?

$\bar{x} \in \mathbb{R}^2$

$\phi(\bar{x}) \in \mathbb{R}^?$

$$K(x, z) = (1 + x_1 z_1 + x_2 z_2)^3$$
$$= \ldots$$
$$= < 1, x_1, x_2, x_1^2, x_2^2, x_1^2 x_2, x_1 x_2^2, x_1^3, x_2^3, x_1 x_2 >$$

10 dimensional?

Q) For $\bar{x} = x$; what space does RBF kernel lie in?

$$K(x,z) = e^{-\gamma||x-z||^2}$$
$$= e^{-\gamma(x-z)^2}$$

Now:

$$e^{\alpha} = \sum_{n=0}^{\infty} \frac{\alpha^n}{n!}$$

$\therefore e^{-\gamma(x-z)^2}$ is $\infty$ dimensional!!

## SVM: Parametric or Non-Parametric

Q) Is SVM parametric or non-parametric?

## SVM: Parametric or Non-Parametric

Q) Is SVM parametric or non-parametric?

Yes and No

Yes $\rightarrow$ Linear kernel or polynomial kernel (form fixed)

No $\rightarrow$ RBF (form changes with data)

## RBF is Non-Parametric

$$\hat{y}(x_{test}) = sign(\bar{w}\bar{x}_{test} + b)$$
$$= sign(\sum_{j=1}^{N_{SV}} \alpha_j y_j \bar{x}_j \bar{x}_{test} + b)$$
$$\hat{y}(X_{test}) = sign(\sum_{j=1}^{N} \alpha_j y_j K(\bar{x}_j, \bar{x}_{test}) + b)$$
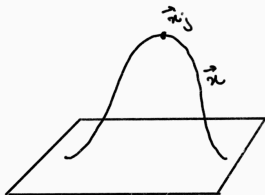
$\alpha_j = 0$ where $j \neq$ S.V.

## RBF is Non-Parametric

Now $K(\bar{x}_j, \bar{x}_{test})$ for RBF is:

$$e^{-\gamma ||\bar{x}_j - \bar{x}_{test}||^2}$$

$\therefore$ Hypothesis is a function of "all" train points

Closer $\bar{x}$ is to $\bar{x}_N$; more is it influencing $\hat{y}(\bar{x})$ - hypothesis function



$\gamma = Low$

High influence of $\bar{x}_j$

- Now if we add a point to the dataset

## RBF is Non-Parametric

- Now if we add a point to the dataset
- Functional form can adapt (similar to KNN)

## RBF is Non-Parametric

- Now if we add a point to the dataset
- Functional form can adapt (similar to KNN)
- $\therefore$ SVM with RBF kernel is non-parametric

## Interpretation of RBF

- $\hat{y}(x) = sign(\sum \alpha_i y_i e^{-||x-x_i||^2} + b)$

## Interpretation of RBF

- $\hat{y}(x) = sign(\sum \alpha_i y_i e^{-||x-x_i||^2} + b)$
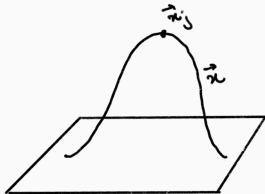- $-||x - x_i||^2$ corresponds to radial term

## Interpretation of RBF

- $\hat{y}(x) = sign(\sum \alpha_i y_i e^{-||x-x_i||^2} + b)$
- $-||x - x_i||^2$ corresponds to radial term
- $\sum \alpha_i y_i$ is the activation component
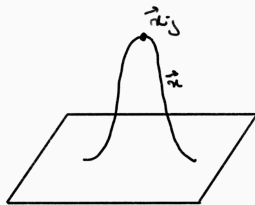
## Interpretation of RBF

- $\hat{y}(x) = sign(\sum \alpha_i y_i e^{-||x-x_i||^2} + b)$
- $-||x - x_i||^2$ corresponds to radial term
- $\sum \alpha_i y_i$ is the activation component
- $e^{-||x-x_i||^2}$ is the basis component

$\gamma$: How far is the influence of a single training sample