# Inductive Bias of SVCs

S1 (band) which is more narrow
distance between Classes is smaller.



margin
S1
(S)
S2
P
6

Support vector Classifiers.

binary, linear
(use line, hyperplane to separate).

Support vector
Decision Boundary → center of the band
Support vector                    (most fair)

S has same slope parallel lines,
Given a slope, if we look at the band nearby

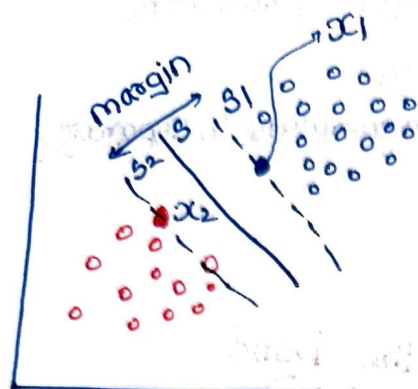Given a slope → we looks at bands → choose the one in
                                      middle.

Choose slope: (that produces) maximum margin
                                  (not narrow)

                    have a little chance of error
                                    (space)

# 2/ mini Review of Basic math.

## Parallel Hyperplanes and their equations



$$x\omega^T + b = 0$$

$$x_1 \omega^T + b = C$$
↳ point not on line.

multiply by $C^{-1}$

$$x_1 \frac{\omega^T}{C} + \frac{b}{C} = 1 \to x\omega^T + b = 0$$

we still get the same point

$$x_1 \omega^T + b = 1$$

$$x_2 \omega^T + b = -1$$

$$x\omega^T + b - 1 = 0 \to S_1$$

$$x\omega^T + b + 1 = \vec{0} \to S_2$$

## Distance between two hyperplanes



$$S_1 : x\omega^T + b - 1 = 0$$

distance $(S_1, 0) \to$ origin $= \dfrac{|b-1|}{\|\omega\|_2} \to$ intercept

$\|\omega\|_2 \to$ 2norm of weights

distance $(S_2, 0) \cdot \dfrac{|b+1|}{\|\omega\|_2} \to$ intercept

$\|\omega\|_2 \to$ norm of weights

Subtract them $\quad d(S_1, S_2) = \dfrac{|b+1 - (b-1)|}{\|\omega\|_2} \cdot \dfrac{2}{\|\omega\|_2}$

## Inner Dot Product as a similarity measure!-

$x_1, x_2 \quad \|x_1 - x_2\|_2^2 = \underset{row}{(x_1-x_2)} \underset{column}{(x_1-x_2)^T}$

euclidean$^2$ $(x_1, x_2)$ ←

inner product of vector

$$= \underset{\underset{1st}{\vee}}{\|x_1\|^2} + \underset{\underset{2nd}{\vee}}{\|x_2\|^2} - \underset{\underset{inner\ product}{\vee}}{2 x_1^T x_2}$$

$(x_1 x_2)^T \longrightarrow$ becomes similarity
dropping -re sign $\qquad \rightsquigarrow$ euclidean similarity

* Similarity measures related to engineered features.

good $\rightarrow$ linear regression $\left. \vphantom{\begin{matrix}a\\b\end{matrix}}\right\}$ non linear classification
$\phantom{good \rightarrow}$ ↳ perceptron

KNN

Similarity → $\quad k(x_1, x_2) = \exp(-\|x_1 - x_2\|_2)$

Similarity Measure

$\phantom{x}$ euclidean similarity

This case. $\qquad k(x_1, x_2) = \varphi(x_1) \cdot \varphi(x_2^T) \rightarrow$ after transformation.

$\qquad\qquad\qquad\qquad$ ↳

$\qquad\qquad\qquad$ adds engineered features.

# 3) The objective in SVCs



$$S_1 \to X\omega^T + b = 1$$
$$S_2 \to S\omega^T + b = -1$$

$$X = [\underbrace{X_1, X_2}_{}]$$

dataset divided in 2 parts

$$\left. \begin{array}{l} x' \in X_1 : x'\omega^T + b \geqslant 1 \end{array} \right\}$$
all points in Blue should be.

for red

$$x' \in X_2 : x'\omega^T + b \leq -1$$

n constraints : Satisfied.

$$\text{margin} = \frac{2}{\|\omega\|_2} \Big\} \text{ maximized.}$$

$$\left. \text{minimize} : \frac{\|\omega\|_2}{2} \right\} \begin{array}{l} \text{In general I hope I get} \\ \text{small weights} \end{array}$$
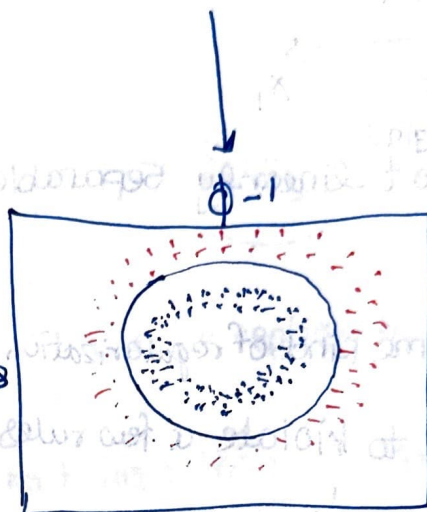
minimize weights
maximize margin.

# Optimization and Kernel Trick



$$x' \in X_1 \rightarrow x'\omega^T + b \geqslant 1$$
$$x^2 \in X_2 \rightarrow x'\omega^T + b \leq -1$$
Linear inequalities

$$\min \|\omega\|_2 \rbrack \text{ quadratic function}$$

Quadratic Programming } not so efficient

## Kernel Matrix

$$n \times n \quad Q_{i,j} = \begin{cases} X_i \cdot X_j^T \rightarrow \text{same label.} \\ -X_i \cdot X_j^T \rightarrow \text{different label} \end{cases} \text{ even less efficient}$$

{ kernel matrix

minimize $v^T Q v$

↓
v (some vector)

→ Solve problem based on similarities (inner product).

**example**

$x_1 = +1 \quad x_2 = +1$
$\begin{cases} K(x_1, x_2) \sim 1 \\ x_1 = +1 \quad x_2 = -1 \\ K(x_1, x_2) \sim 0 \end{cases}$ Similarity

Kernel Trick.
I can find my SVC → euclidean similary
redefine the definition of similarity

SVC classifier becomes _____

$$K_{i,j} = \exp\left(-\|X_i - X_j\|_2\right)$$

→ non linear

# Slacks



outlier
(forcing margin)

not a
good thing

not linearly separable → SVC fail

↓

introduce some kind of regularization

Sensitivity

allow points to violate a few rules

↓

introduce slacks

$X: X_{pos}, X_{neg}$

TOTAL SLACK

$$\xi \sum_{i=1}^{n} \xi_i$$

$$\min \|w_2\|^2 + C \cdot \xi$$

if C is small
   give more slack

if C big don't give
      too much
      slack.

↓
strength of
regularization

$X_i \in X_{pos}: X_1 w^T + b > 1 - \overline{\xi_i}$

introduce $C_i$ for each $X_i$

↓ $\overline{C_i} > 0$

+ve number

so you don't force it on one side
         you allow it to be on wrong side

$\xi_i$ for each $X_i \in X_{neg}$

↓ $X_i \in X_{neg}: Y_i w^T + b \leq -1 + \overline{\xi_i}$

if $\overline{\xi_i}$ is greater point is allowed to even be on other
      than 1                side, so its fine

thats how much slack we
give.

non differentiable
because of other variations.

called hinge loss

$$\xi_i$$

$$\min \|w\|_2^2 + C \left( \sum_{i=1}^{n} \max \left(0, 1 - \text{sign}(y_i - 0.5)(xw^T + b)\right) \right)$$

$\leftarrow$ equivalent

$$\min \|w\|_2^2 + C \left( \sum_{i=1}^{n} \max \right)$$

$$\min \|w\|_2^2 + C \sum_i \xi_i \left(\xi_i > 0\right) \bigg] \text{loss}.$$

$\underbrace{\qquad\qquad}$

Loss function

$\downarrow$

can be replaced by

when we have Slack, we are not just having regularization

we are also doing gradient descent.

$\downarrow$

when we have linear kernel.

$\Bigg]$

faster linear
SVC

# Hinge Loss

non differentiable
because of other
variations.

called hinge loss

$$\min \|w\|_2^2 + C\left(\sum_{i=1}^{n} \max\left(0, 1 - \text{sign}(y_i - 0.5)(xw^T + b)\right)\right)$$

$\xi_i$

← equivalent

$$\min \|w\|_2^2 + C\left(\sum_{i=1}^{n} \max \right.$$

$$\min \|w\|_2^2 + C\sum_i \xi_i \ (\xi_i > 0)] \ \text{loss}.$$

$$\underbrace{\phantom{C\sum_i \xi_i}}_{\text{Loss function}}$$

can be replaced by

when we have Slacks, we are not just having regularization
we are also doing gradient descent.

when we have linear kernel.

faster linear
SVC