

Support Vector Machine (SVM)

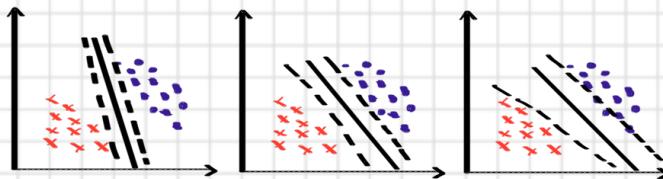
- IT IS A **SUPERVISED LEARNING METHOD**

- CAN BE USED BOTH FOR:

- **CATEGORIZATION**
- **REGRESSION**

- INTUITION:

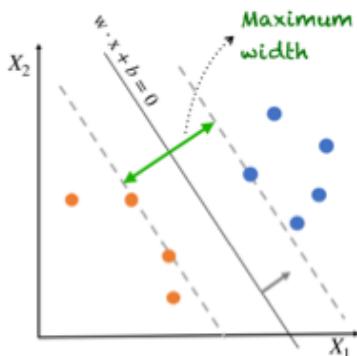
- WHAT'S THE **BEST DECISION BOUNDARY** FOR THESE TRAINING DATA POINTS?



- IF CLASSES ARE **LINEARLY SEPARABLE**:

* BEST DECISION BOUNDARY \rightarrow LINE (**HYPERSPACE**)

* THE ONE WITH **MAXIMUM MARGIN** FROM CLOSEST SAMPLES OF BOTH CLASSES!



WE BEGIN HERE: **HARD MARGIN SVM**

THE "STREET" THAT FORMS INSIDE IS
CALLED **WIDEST STREET** AND THE NAME
INDUCES THE **UNDERLINED** MEANING

- ORIGINALLY SVM IS USED FOR **BINARY CLASSIFICATION**

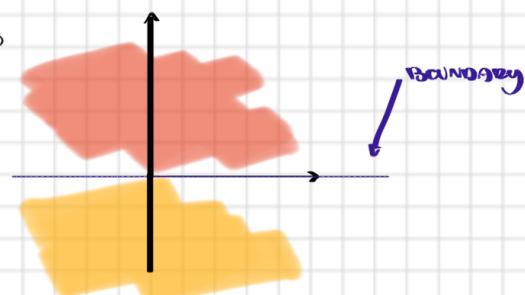
- COMBINING SEVERAL SVM'S FOR **MULTI-CLASS** IS POSSIBLE THOUGH

- **LINEAR CLASSIFIER**:

$$y(x) = w \cdot x + b \quad , \quad \text{DECISION BOUNDARY} : y(x) = 0$$

IF $y(x_i) \geq 0$ THEN SAMPLE i IS IN CLASS **+1**

IF $y(x_i) < 0$ THEN SAMPLE i IS IN CLASS **-1**



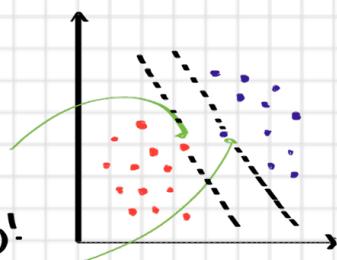
Hard Margin

- THINK ABOUT WHICH TRAINING SAMPLES ARE RELEVANT FOR **WIDEST STREET**?

* THE TRAINING SAMPLES THAT ARE RELEVANT ARE THE ONLY ONES THAT CONTRIBUTE TO THE DECISION BOUNDARY.

↳ **Support Vectors!**

* SVM **FINDS** SUPPORT VECTORS DURING TRAINING & USES THEM TO CLASSIFY!



* LINEAR CLASSIFIER:

IF OF TRAINING SAMPLES!

$$y(x) = w \cdot x + b$$

FEATURES VECTOR & OUTPUT OF TRAINING SAMPLE i

$$y(x) = \sum_{i=1}^n w_i y_i (x_i \cdot x) + b$$

$w_i \geq 0$ input vector to be classified!

dot product between training datapoints & new datapoint

Some samples
may HAVE
 $y_i = 0$.
THE ONES WITH
 $y_i > 0$ ARE SUPPORT VECTORS

TRAINING AN SVM MEANS FINDING VALUES FOR ALL THE w_i 'S, EACH ASSOCIATED WITH A TRAINING SAMPLE SUCH THAT THE WIDTH OF THE STREET SEPARATING THE TWO CLASSES IS MAXIMIZED.

THIS MEANS FIGURING OUT WHICH $y_i = 0$, IDENTIFYING SUP. VECTORS.

AFTER TRAINING KEEP ONLY SUP. VEC. ASSOCIATED WITH y_i FOR PREDICTION

• So → TRAIN THE MODEL = FIND THE w_i MULTIPLIERS ✓ SINGLE TRAINING SAMPLE

THIS LEADS TO = MAXIMIZE STREET WIDTH

• KEEP ONLY SUPPORT VECTORS (TRAIN SAMPLES WITH NON-ZERO w_i) & THEIR y_i

• USE SUP. VEC. & THEIR w_i TO CLASSIFY NEW SAMPLES.

* WE USE HARD MARGIN WHEN TWO CLASSES ARE LINEARLY SEPARABLE

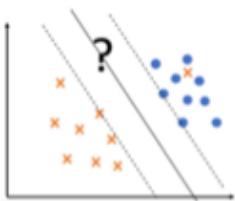
* WE FIND w_i 'S SO THAT WIDTH OF // IS MAXIMIZED

* NO TRAINING SAMPLES ARE ALLOWED ON THE WRONG SIDE OF MARGINS

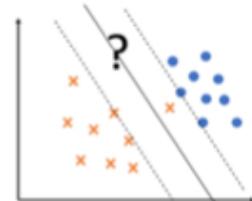
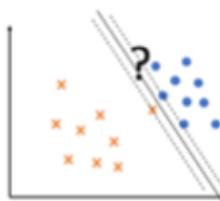
From here now a problem arises:

But what if the class conditional distributions overlap?

Or if they look like this?



Impossible to fit a hard-margin SVM for this dataset



SOFT MARGINS

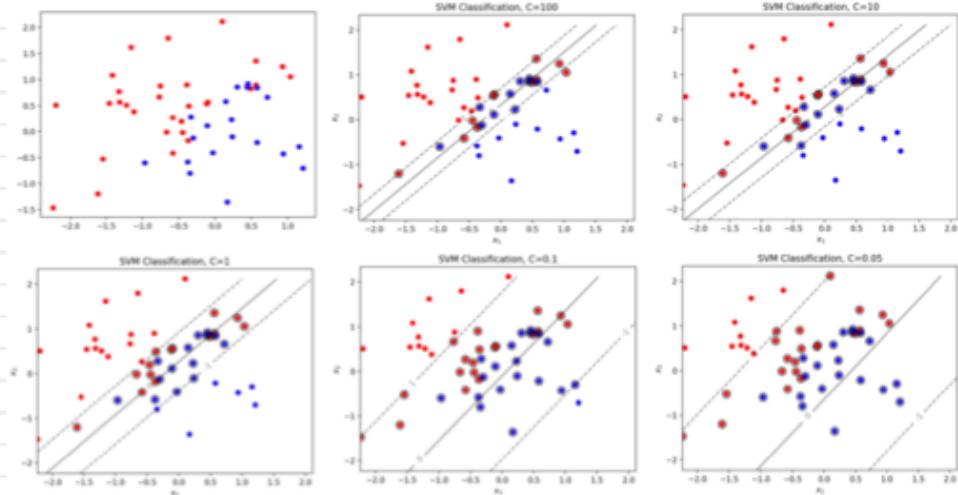
- CONSTRAINT IS RELAXED SO IT IS ALLOWED FOR SOME OF THE TRAINING SAMPLES TO FALL ON WRONG SIDE.
- WE DO AWAY IT BUT WE PENALIZE IT!!
- OPTIMIZATION = MAXIMIZE STREET WHILE MINIMIZING PENALTY

ONTO THE MATH:

$$y(x) = \sum_{i=1}^n w_i y_i (x_i \cdot x) + b \quad \text{OR} \quad C$$

THIS LOOKS THE SAME AS HARD MARGIN EXCEPT THE CONDITION FOR w_i 'S.

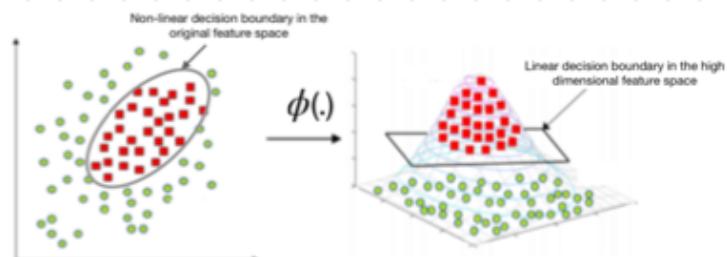
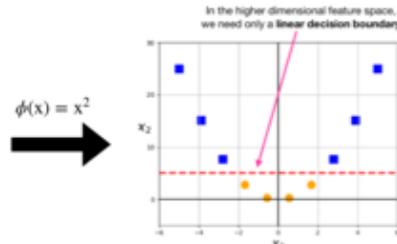
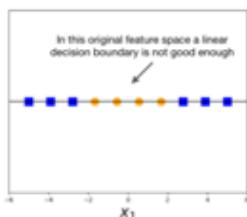
- THE HYPERPARAMETER $C > 0$ MAKES A TRADEOFF BETWEEN THE PENALTY FOR THE SAMPLES ON THE WRONG SIDE OF THE MARGINS & MAX WIDTH
- ALL SAMPLES ON THE MARGINS OR ON THE WRONG SIDE OF THE MARGINS CONTRIBUTE



SVM WITH NON-LINEAR DECISION BOUNDARIES

- SOMETIMES THE SPACE MIGHT BE NON-LINEAR:

- MAPPING THE FEATURES FROM ORIGINAL DM TO HIGHER ONES
- THIS MAKES IT LINEAR



A LINEAR DECISION BOUNDARY IN THE HIGH DIM FEATURE SPACE CAN REPRESENT A NON-LINEAR DECISION BOUNDARY IN ORIGINAL SPACE

- SOME TRANSFORMATION CAN ADD MANY NEW NON-LINEAR FEATURES

↙ THIS IS DUE TO HOPEFULLY FIND A GOOD LINEAR DECISION BOUNDARY

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \rightarrow \phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1 x_2 \\ x_1 x_2 \dots x_d \\ \vdots \end{bmatrix}$$

USING TRANSFORMATION $\phi(\cdot)$ WE MAP THE FEATURES FROM THEIR LOW DIMENSIONAL SPACE TO A HIGH DIMENSIONAL SPACE WHERE PROBS IS LINEAR.

MODEL:

$$y(x) = \sum_i^n y_i (\phi(x_i) \cdot \phi(x)) + b$$

- BUT A VERY HIGH DIMENSION MEANS TOO MANY COMPUTATIONS REQUIRED FOR DOT PRODUCTS
- $\phi(x_i) \cdot \phi(x)$ DOT PRODUCT HAS A TIME OF ORDER $O(2^d)$

KERNELS

- MAKES NON-LINEAR DECISION BOUNDARIES EASY TO CREATE

$$y(x) = \sum_i^n y_i K(x_i, x) + b$$

KERNEL TRICK

For some transformations $\phi(x)$ there is a kernel function $K(x, z)$ that can take any pair of input vectors & compute the dot product of their transformed version ($\phi(x) \cdot \phi(z)$) directly in the original feature space

We could transform feature vectors to infinite dims & still have no extra computation

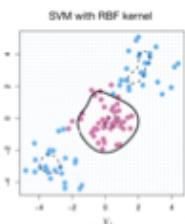
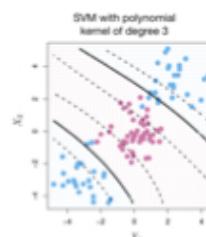
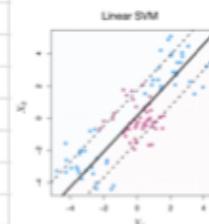
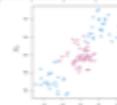
We just need to be able to compute $K(x, z)$

Popular ones:

- Linear Kernel: $K(x, z) = x \cdot z$

Polynomial Kernel: $K(x, z) = (x \cdot z + 1)^b$

Example:



Pros vs Cons

Pros

- ^ EFFECTIVE IN HIGH DIMENSIONAL SPACE
- ^ " WHEN # OF FEATURES > TRAINING EXAMPLES
- ^ BEST ALGO WHEN CLASSES ARE SEPARABLE
- ^ MEMORY EFFICIENT DUE TO SUPPORT VECTORS
- ^ OUTLIERS HAVE LESS IMPACT "
- ^ RESULTS CAN BE EXPLAINED!!

Cons

- ^ LARGE DATASET - TAKES A LOT TO TRAIN!
- ^ SENSITIVE TO HYPER PARAMS
- ^ SELECTING CORRECT/APPROPRIATE KERNEL