

Lecture 17: 31 March, 2022

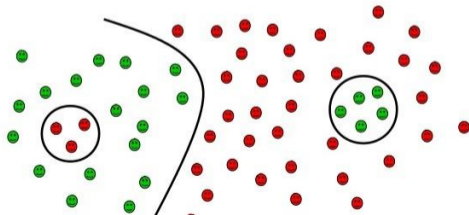
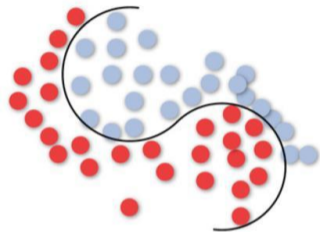
Madhavan Mukund

<https://www.cmi.ac.in/~madhavan>

Data Mining and Machine Learning
January–May 2022

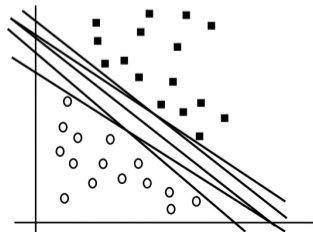
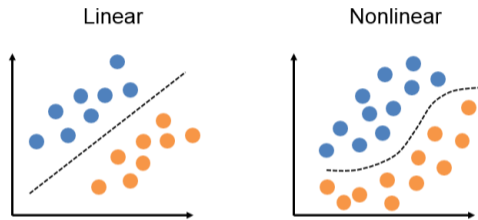
A geometric view of supervised learning

- Think of data as points in space
- Find a separating curve (surface)
- Separable case
 - Each class is a connected region
 - A single curve can separate them
- More complex scenario
 - Classes form multiple connected regions
 - Need multiple separators



Linear separators

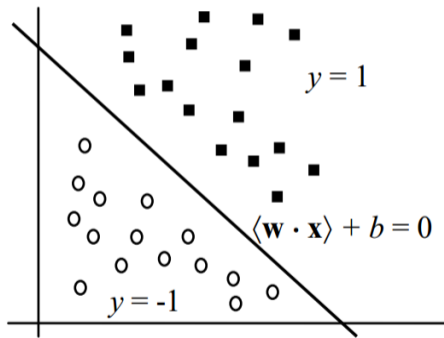
- Simplest case – linearly separable data
- Dual of linear regression
 - Find a line that passes close to a set of points
 - Find a line that separates the two sets of points
- Many lines are possible
 - How do we find the best one?
 - What is a good notion of "cost" to optimize?



Linear separators

- Each input x has n attributes $\langle x_1, x_2, \dots, x_n \rangle$
- Linear separator has the form
$$w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$
- Classification criterion
$$w_1x_1 + \dots + w_nx_n + b > 0, \text{ classify yes, } +1$$
$$w_1x_1 + \dots + w_nx_n + b < 0, \text{ classify no, } -1$$
- Dot product $\langle w \cdot x \rangle$
$$(w_1, \dots, w_n) \cdot (x_1, \dots, x_n) = w_1x_1 + \dots + w_nx_n$$
- Collapsed form
$$\langle w \cdot x \rangle + b > 0, \langle w \cdot x \rangle + b < 0$$
- Rename bias b as w_0 , create fictitious $x_0 = 1$
- Equation becomes

$$\langle w \cdot x \rangle > 0, \langle w \cdot x \rangle < 0$$



Perceptron algorithm

(Frank Rosenblatt, 1958)

- Each training input is (x_i, y_i) where $x_i = \langle x_1^i, x_2^i, \dots, x_n^i \rangle$ and $y_i = +1$ or -1
- Need to find $w = \langle w_0, w_1, \dots, w_n \rangle$.
Recall $x_0^i = 1$, always

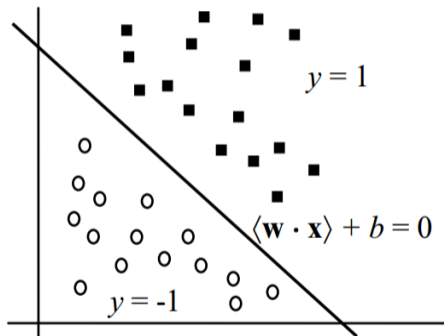
Initialize $w = \langle 0, 0, \dots, 0 \rangle$

While there exists (x_i, y_i) such that

$y_i = +1$, and $\langle w \cdot x_i \rangle < 0$, or

$y_i = -1$, and $\langle w \cdot x_i \rangle > 0$

Update w to $w + x_i y_i$



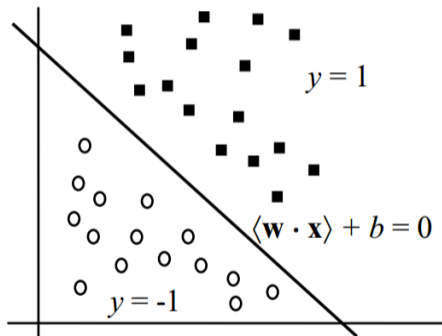
Perceptron algorithm

- Keep updating w as long as some training data item is misclassified
- Update is an offset by misclassified input
- Need not stabilize, potentially an infinite loop

Theorem

If the points are linearly separable, the Perceptron algorithms always terminates with a valid separator

- Termination time depends on two factors
 - Width of the band separating the positive and negative points
 - Narrow band takes longer to converge
 - Magnitude of the x values
 - Larger spread of points takes longer to converge



Perceptron Algorithm — Proof

Theorem

If there is w^* satisfying $(w^* \cdot x_i)y_i \geq 1$ for all i , then the Perceptron Algorithm finds a solution w with $(w \cdot x_i)y_i > 0$ for all i in at most $r^2|w^*|^2$ updates, where $r = \max_i |x_i|$.

$$\begin{aligned}w \cdot x_i &> 0 & y_i = 1 \\w \cdot x_i y_i &> 0 \\w \cdot x_i < 0 & w \cdot x_i y_i > 0 & y_i = -1\end{aligned}$$

$$\begin{aligned}\lambda w \cdot x_i &> 0 \text{ if } y_i = 1 \\ \lambda w \cdot x_i &< -1 \text{ if } y_i = -1\end{aligned}$$

Perceptron Algorithm — Proof

Theorem

If there is w^* satisfying $(w^* \cdot x_i)y_i \geq 1$ for all i , then the Perceptron Algorithm finds a solution w with $(w \cdot x_i)y_i > 0$ for all i in at most $r^2|w^*|^2$ updates, where $r = \max_i |x_i|$.

- Assume w^* exists. Keep track of two quantities: $w^T w^*$, $|w|^2$.

fixed
our incrementally updated estimate

Perceptron Algorithm — Proof

Theorem

If there is w^* satisfying $(w^* \cdot x_i)y_i \geq 1$ for all i , then the Perceptron Algorithm finds a solution w with $(w \cdot x_i)y_i > 0$ for all i in at most $r^2|w^*|^2$ updates, where $r = \max_i |x_i|$.

- Assume w^* exists. Keep track of two quantities: $w^T w^*$ and $|w|^2$.
- Each update increases $w^T w^*$ by at least 1.

$$\underbrace{(w + x_i y_i)^T}_{\text{new } w} w^* = \underbrace{w^T}_{\text{green}} w^* + \underbrace{x_i^T y_i w^*}_{\geq 1} \geq \underbrace{w^T w^*}_{\text{red}} + \underbrace{1}_{\text{red}}$$

Perceptron Algorithm — Proof

Theorem

If there is w^* satisfying $(w^* \cdot x_i)y_i \geq 1$ for all i , then the Perceptron Algorithm finds a solution w with $(w \cdot x_i)y_i > 0$ for all i in at most $r^2|w^*|^2$ updates, where $r = \max_i |x_i|$.

- Assume w^* exists. Keep track of two quantities: $w^T w^*$, $|w|^2$.

- Each update increases $w^T w^*$ by at least 1.

$$(w + x_i y_i)^T w^* = w^T w^* + x_i^T y_i w^* \geq w^T w^* + 1$$

- Each update increases $|w|^2$ by at most r^2 .

$$(w + x_i y_i)^T (w + x_i y_i) = |w|^2 + 2x_i^T y_i w + |x_i y_i|^2 \leq |w|^2 + |x_i|^2 \leq |w|^2 + r^2$$

- Note that we update only when $x_i^T y_i w < 0$

$$y_i = 1 \quad w \cdot x_i < 0$$

$$y_i = -1 \quad w \cdot x_i > 0$$

Perceptron Algorithm — Proof (cont'd)

- Assume Perceptron Algorithm makes m updates

Perceptron Algorithm — Proof (cont'd)

- Assume Perceptron Algorithm makes m updates
- Then, $w^T w^* \geq m$, $|w|^2 \leq mr^2$

↑
increased by
at least 1
m times

↘
increased
by at least r^2
m times

Initially $w = 0$
 $w^T w^* = 0$
 $|w|^2 = 0$

Perceptron Algorithm — Proof (cont'd)

- Assume Perceptron Algorithm makes m updates
- Then, $w^\top w^* \geq m$, $|w|^2 \leq mr^2$
- $m \leq |w||w^*|$

Perceptron Algorithm — Proof (cont'd)

- Assume Perceptron Algorithm makes m updates
- Then, $w^\top w^* \geq m$, $|w|^2 \leq mr^2$
- $$m \leq |w||w^*|$$
$$m/|w^*| \leq |w|$$

Perceptron Algorithm — Proof (cont'd)

- Assume Perceptron Algorithm makes m updates

- Then, $w^\top w^* \geq m$, $|w|^2 \leq mr^2$

- $m \leq |w||w^*|$

$$m/|w^*| \leq |w|$$

$$m/|w^*| \leq r\sqrt{m}$$

Perceptron Algorithm — Proof (cont'd)

- Assume Perceptron Algorithm makes m updates

- Then, $w^\top w^* \geq m$, $|w|^2 \leq mr^2$

- $m \leq |w||w^*|$

$$m/|w^*| \leq |w|$$

$$m/|w^*| \leq r\sqrt{m}$$

$$\sqrt{m} \leq r|w^*|$$

Perceptron Algorithm — Proof (cont'd)

- Assume Perceptron Algorithm makes m updates

- Then, $w^\top w^* \geq m$, $|w|^2 \leq mr^2$

- $m \leq |w||w^*|$

$$m/|w^*| \leq |w|$$

$$m/|w^*| \leq r\sqrt{m}$$

$$\sqrt{m} \leq r|w^*|$$

$$m \leq r^2|w^*|^2$$

magnitude of points
inverse of width of separating region

Perceptron Algorithm — Proof (cont'd)

- Assume Perceptron Algorithm makes m updates

- Then, $w^\top w^* \geq m$, $|w|^2 \leq mr^2$

- $m \leq |w||w^*|$

$$m/|w^*| \leq |w|$$

$$m/|w^*| \leq r\sqrt{m}$$

$$\sqrt{m} \leq r|w^*|$$

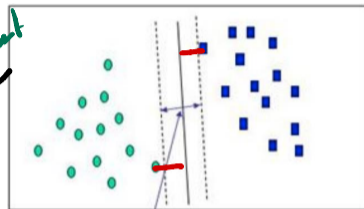
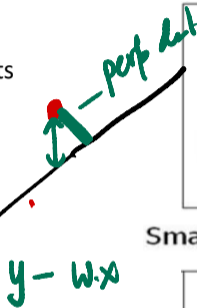
$$m \leq r^2|w^*|^2$$

- Note (for later) that final w is of the form $\sum_i n_i x_i$

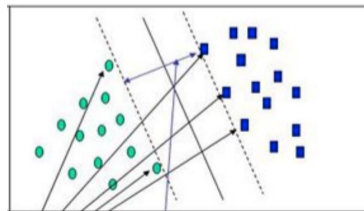
<u>Iteration</u>	<u>w</u>
0	0
1	$\pm x_{i_1}$
2	$\pm x_{i_1} \pm x_{i_2}$
	$x_{i_1} - x_{i_2} + x_{i_3} - x_{i_4}$

Margin

- Each separator defines a *margin*
 - Empty corridor separating the points
 - Separator is the centre line of the margin
- Wider margin makes for a more robust classifier
 - More gap between the classes
- Optimum classifier is one that maximizes the width of its margin
- Margin is defined by the training data points on the boundary
 - Support vectors



Small Margin



Large Margin

Support Vectors

Finding a maximum margin classifier

- Recall our original linear classifier

$$w_1x_1 + \dots + w_nx_n + b > 0, \quad \text{classify yes, } +1$$

$$w_1x_1 + \dots + w_nx_n + b < 0, \quad \text{classify no, } -1$$

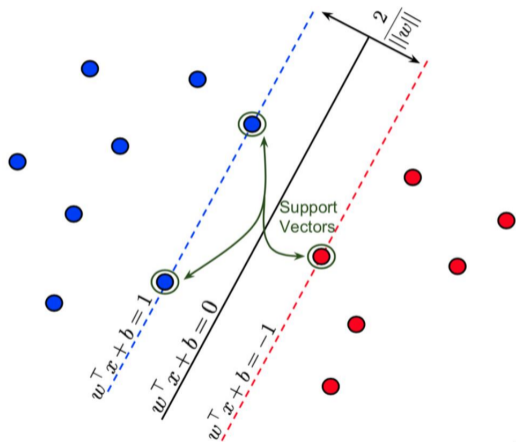
- Scale margin so that separation is 1 on either side

$$w_1x_1 + \dots + w_nx_n + b > 1, \quad \text{classify yes, } +1$$

$$w_1x_1 + \dots + w_nx_n + b < -1, \quad \text{classify no, } -1$$

- Using Pythagoras's theorem, perpendicular distance to nearest support vector is $\frac{1}{\|w\|}$,

$$\text{where } \|w\| = \sqrt{w_1^2 + w_2^2 + \dots + w_n^2}$$

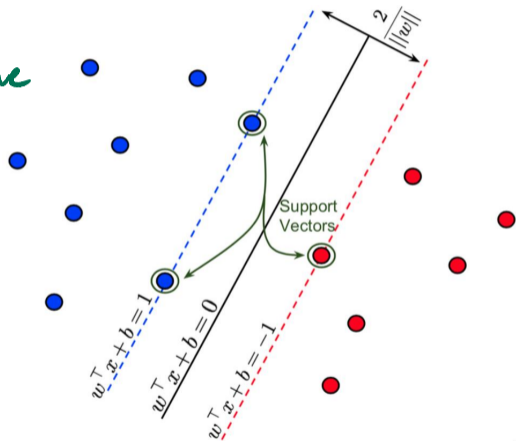


Optimization problem

- Want to maximize the overall margin $\frac{2}{\|w\|}$
- Equivalently, minimize $\frac{\|w\|}{2}$ — Objective
- Also, w should classify each (x_i, y_i) correctly

$$\begin{cases} w_1 x_1^i + \dots + w_n x_n^i + b > 1, & \text{if } y_i = 1 \\ w_1 x_1^i + \dots + w_n x_n^i + b < -1, & \text{if } y_i = -1 \end{cases}$$

constraints



Optimization problem

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

Subject to

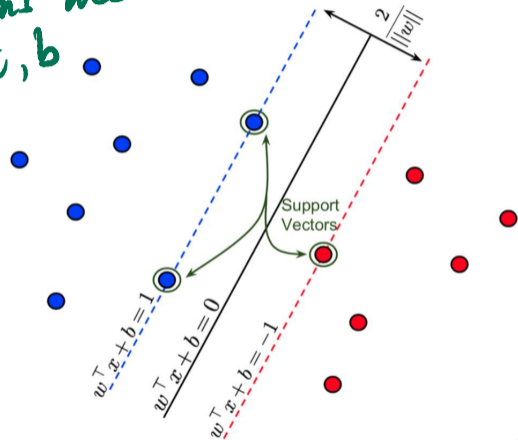
$w_1 x_1^i + \dots + w_n x_n^i + b > 1, \quad \text{if } y_i = 1$
 $w_1 x_1^i + \dots + w_n x_n^i + b < -1, \quad \text{if } y_i = -1$

- The objective function is not linear

$\|w\| = \sqrt{w_1^2 + w_2^2 + \dots + w_n^2}$

- This is a *quadratic optimization* problem, not linear programming

Unknowns are w_i, b



Linear constraints



Solution to optimization problem

input (x_1, y_1) to (x_N, y_N)

- Convex optimization theory
- Can be solved using computational techniques
- Solution expressed in terms of Lagrange multipliers $\alpha_1, \alpha_2, \dots, \alpha_N$
one multiplier per training input
- α_i is non-zero iff x_i is a support vector
- Final classifier for new input z

$$\text{sign} \left[\sum_{i \in sv} y_i \alpha_i \langle x_i \cdot z \rangle + b \right]$$

- sv is set of support vectors

α_1 Constraint 1

α_2 Constraint 2

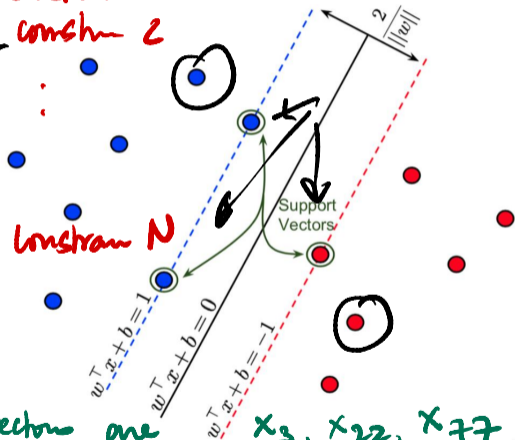
⋮

α_N Constraint N

Support vectors one

x_3, x_{22}, x_{77}

$$y_3 \alpha_3 \langle x_3 \cdot z \rangle + y_{22} \alpha_{22} \langle x_{22} \cdot z \rangle + y_{77} \alpha_{77} \langle x_{77} \cdot z \rangle + b$$

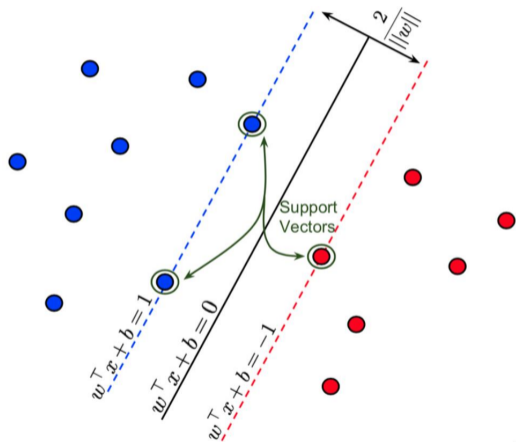


Support Vector Machine (SVM)

$$\text{sign} \left[\sum_{i \in sv} y_i \alpha_i \langle x_i \cdot z \rangle + b \right]$$

Support Vector Machine (SVM)

- Solution depends only on support vectors
 - If we add more training data away from support vectors, separator does not change
- Solution uses dot product of support vectors with new point
 - Will be used later, in the non-linear case



The non-linear case

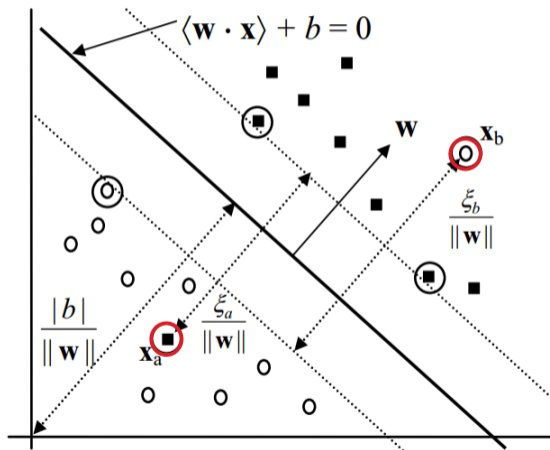
- Some points may lie on the wrong side of the classifier
- How do we account for these?
- Add an error term to the classifier requirement

• Instead of

$$\langle w \cdot x \rangle + b > 1, \quad \text{if } y_i = 1$$
$$\langle w \cdot x \rangle + b < -1, \quad \text{if } y_i = -1$$

we have

$$\langle w \cdot x \rangle + b > 1 - \xi_i, \quad \text{if } y_i = 1$$
$$\langle w \cdot x \rangle + b < -1 + \xi_i, \quad \text{if } y_i = -1$$



Too strict

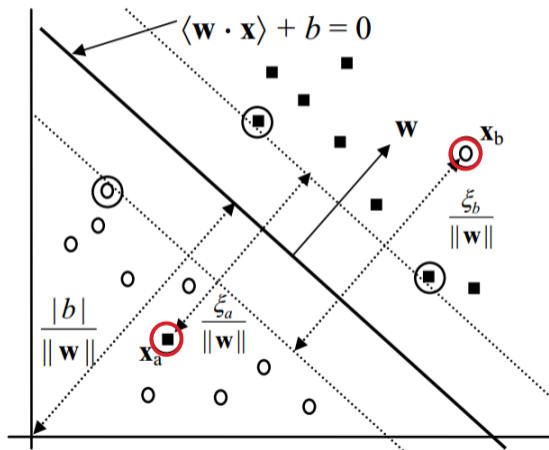
Too generous



Soft margin classifier

$$\langle w \cdot x \rangle + b > 1 - \xi_i, \quad \text{if } y_i = 1$$
$$\langle w \cdot x \rangle + b < -1 + \xi_i, \quad \text{if } y_i = -1$$

- Error term always non-negative, $\xi_i \geq 0$
- If the point is correctly classified, error term is 0
- Soft margin – some points can drift across the boundary
- Need to account for the errors in the objective function
 - Minimize the need for non-zero error terms



Soft margin optimization

$$\text{Minimize } \frac{\|w\|}{2} + \sum_{i=1}^N \xi_i^2$$

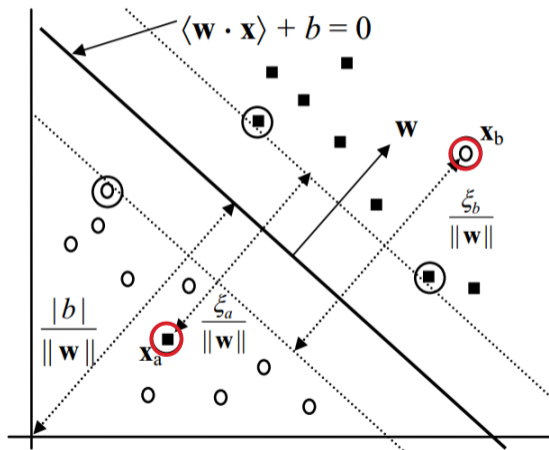
Subject to

$$\xi_i \geq 0$$

$$\langle w \cdot x \rangle + b > 1 - \xi_i, \quad \text{if } y_i = 1$$

$$\langle w \cdot x \rangle + b < -1 + \xi_i, \quad \text{if } y_i = -1$$

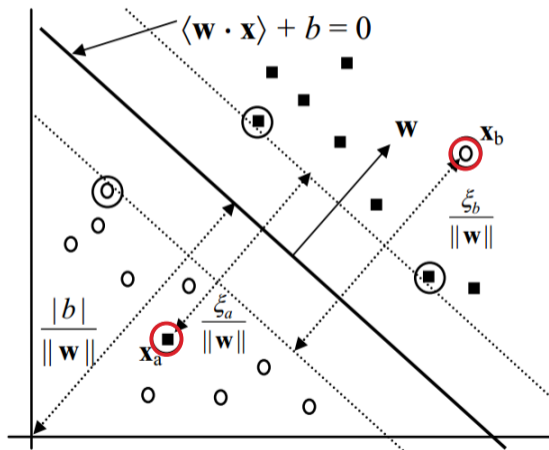
- Constraints include requirement that error terms are non-negative
- Again the objective function is quadratic



Soft margin optimization

- Can again be solved using convex optimization theory
- Form of the solution turns out to be the same as the hard margin case
 - Expression in terms of Lagrange multipliers α_i
 - Only terms corresponding to support vectors are actively used

$$\text{sign} \left[\sum_{i \in sv} y_i \alpha_i \langle x_i \cdot z \rangle + b \right]$$



The non-linear case

- How do we deal with datasets where the separator is a complex shape?
- Geometrically transform the data
 - Typically, add dimensions
- For instance, if we can "lift" one class, we can find a planar separator between levels

