

- The present form of support vector machine (SVM) was largely developed at AT&T Bell Laboratories by Vapnik and co-workers.
- Known as a **maximum margin classifier**.
- Originally proposed for classification and soon applied to regression and time series prediction.
- One of the most efficient **supervised learning** methods.

Problem

- Given a set of training samples

$$(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N), x_i \in R^n, y_i \in \{-1, 1\},$$

find a function $f(x, \alpha)$ to classify the samples, such that

$$f(x_i, \alpha) \begin{cases} > 0, & \forall y_i = +1; \\ < 0, & \forall y_i = -1, \end{cases}$$

where α denotes the parameters.

Problem

- Given a set of training samples

$$(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N), x_i \in R^n, y_i \in \{-1, 1\},$$

find a function $f(x, \alpha)$ to classify the samples, such that

$$f(x_i, \alpha) \begin{cases} > 0, & \forall y_i = +1; \\ < 0, & \forall y_i = -1, \end{cases}$$

where α denotes the parameters.

- For a testing sample x , we can predict its label by $\text{sign}[f(x, \alpha)]$.

Problem

- Given a set of training samples

$$(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N), x_i \in R^n, y_i \in \{-1, 1\},$$

find a function $f(x, \alpha)$ to classify the samples, such that

$$f(x_i, \alpha) \begin{cases} > 0, & \forall y_i = +1; \\ < 0, & \forall y_i = -1, \end{cases}$$

where α denotes the parameters.

- For a testing sample x , we can predict its label by $\text{sign}[f(x, \alpha)]$.
- $f(x, \alpha) = 0$ is called the separation hyperplane.

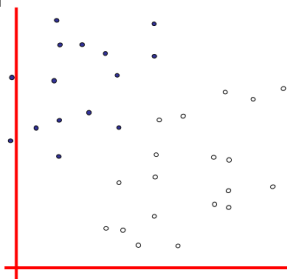
Linear classifiers

Linear hyperplane

$$f(x, w, b) = \langle x, w \rangle + b = 0$$

Consider the linearly separable case, there are infinite number of hyperplanes that can do the job.

- denotes +1
- denotes -1



How would you classify this data?

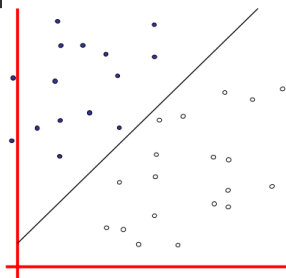
Linear classifiers

Linear hyperplane

$$f(x, w, b) = \langle x, w \rangle + b = 0$$

Consider the linearly separable case, there are infinite number of hyperplanes that can do the job.

- denotes +1
- denotes -1



How would you classify this data?

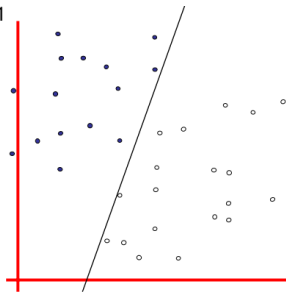
Linear classifiers

Linear hyperplane

$$f(x, w, b) = \langle x, w \rangle + b = 0$$

Consider the linearly separable case, there are infinite number of hyperplanes that can do the job.

- denotes +1
- denotes -1



How would you classify this data?

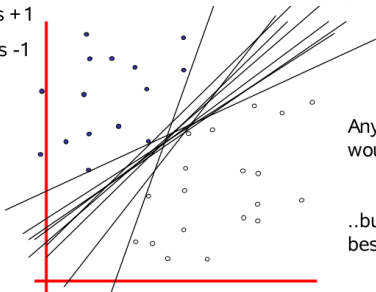
Linear classifiers

Linear hyperplane

$$f(x, w, b) = \langle x, w \rangle + b = 0$$

Consider the linearly separable case, there are infinite number of hyperplanes that can do the job.

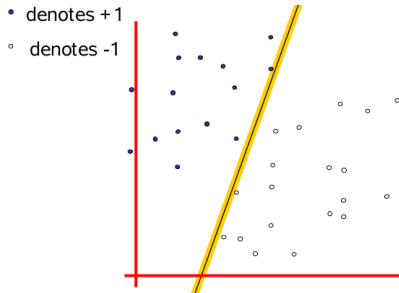
- denotes +1
- denotes -1



Any of these
would be fine..

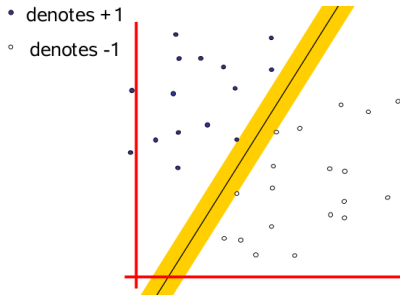
..but which is
best?

Margin of a linear classifier



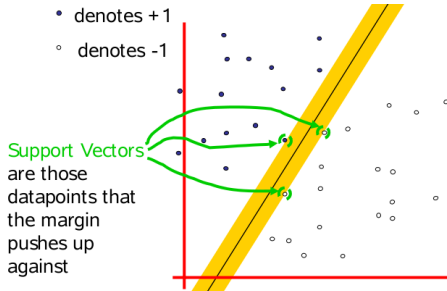
Definition: the width that the boundary could be increased by before hitting a data point.

Maximum margin linear classifier



Definition: the linear classifier with the maximum margin.

Support vectors



Problem formulation

To formulate the margin, we further requires that for all samples

$$f(x_i, \alpha) = \langle x_i, w \rangle + b \begin{cases} \geq +1, & \forall y_i = +1; \\ \leq -1, & \forall y_i = -1. \end{cases}$$

or

$$y_i(\langle x_i, w \rangle + b) \geq 1, \quad i = 1, \dots, N.$$

Problem formulation

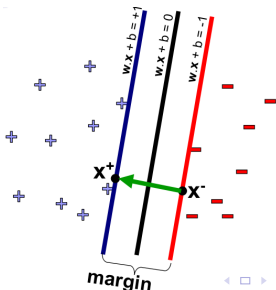
To formulate the margin, we further requires that for all samples

$$f(x_i, \alpha) = \langle x_i, w \rangle + b \begin{cases} \geq +1, & \forall y_i = +1; \\ \leq -1, & \forall y_i = -1. \end{cases}$$

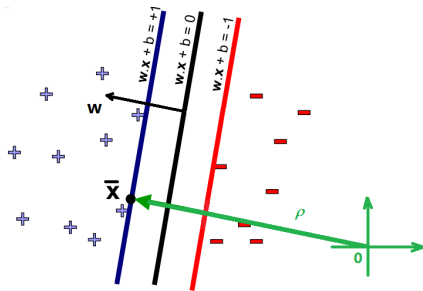
or

$$y_i(\langle x_i, w \rangle + b) \geq 1, \quad i = 1, \dots, N.$$

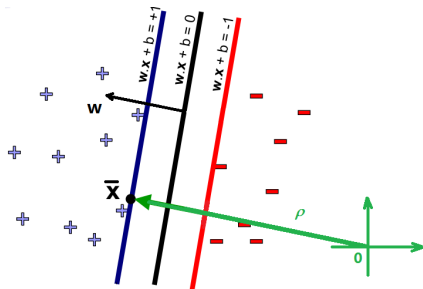
- We have introduced two additional hyperplanes $\langle x, w \rangle + b = \pm 1$ parallel to the separation hyperplane $\langle x, w \rangle + b = 0$



What is the margin? The distance between the two new hyperplanes.

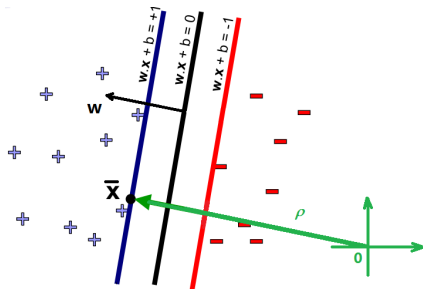


What is the margin? The distance between the two new hyperplanes.



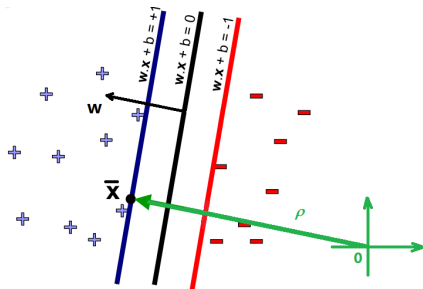
- The minimum distance between the hyperplane $\langle x, w \rangle + b = 1$ and the origin is $\rho_1 = \frac{1-b}{\|w\|}$. (why?)

What is the margin? The distance between the two new hyperplanes.



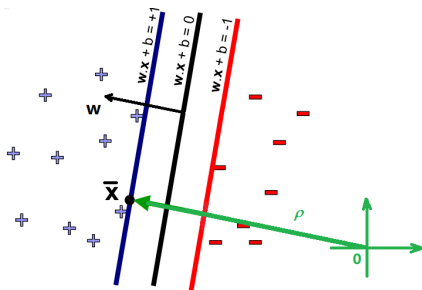
- The minimum distance between the hyperplane $\langle x, w \rangle + b = 1$ and the origin is $\rho_1 = \frac{1-b}{\|w\|}$. (why?)
- The minimum distance between the hyperplane $\langle x, w \rangle + b = -1$ and the origin is $\rho_2 = \frac{-1-b}{\|w\|}$.

What is the margin? The distance between the two new hyperplanes.



- The minimum distance between the hyperplane $\langle x, w \rangle + b = 1$ and the origin is $\rho_1 = \frac{1-b}{\|w\|}$. (why?)
- The minimum distance between the hyperplane $\langle x, w \rangle + b = -1$ and the origin is $\rho_2 = \frac{-1-b}{\|w\|}$.
- The margin is $|\rho_1 - \rho_2| = 2/\|w\|$.

How to calculate ρ_1 and ρ_2 ?

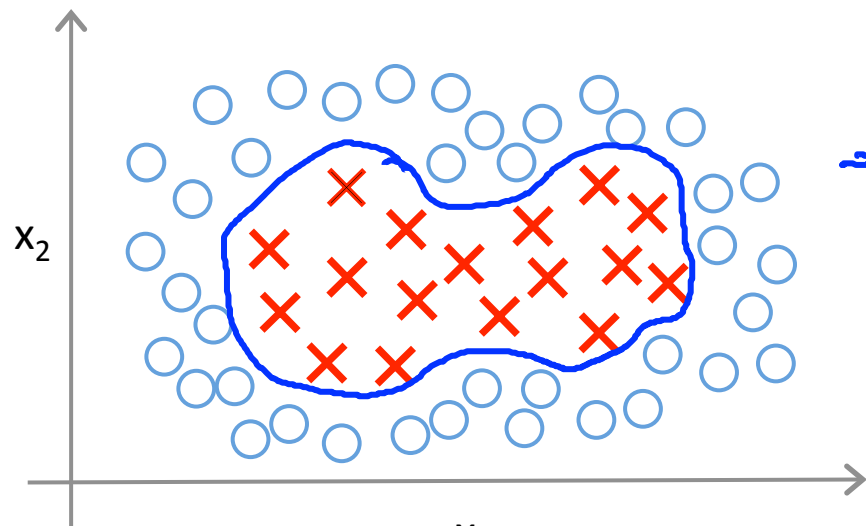


Note $\bar{x} = \rho_1 w / \|w\|$, where $w / \|w\|$ is the unit vector along the direction w . Since \bar{x} is on the blue hyperplane, then

$$\langle \rho_1 w / \|w\|, w \rangle + b = 1$$

which follows $\rho_1 = \frac{1-b}{\|w\|}$. Similarly, we obtain $\rho_2 = \frac{-1-b}{\|w\|}$.

Non-linear Decision Boundary



Predict $y = 1$ if

$$\rightarrow \theta_0 + \theta_1 \underline{x_1} + \theta_2 \underline{x_2} + \theta_3 \underline{x_1 x_2} \\ + \theta_4 \underline{x_1^2} + \theta_5 \underline{x_2^2} + \dots \geq 0$$

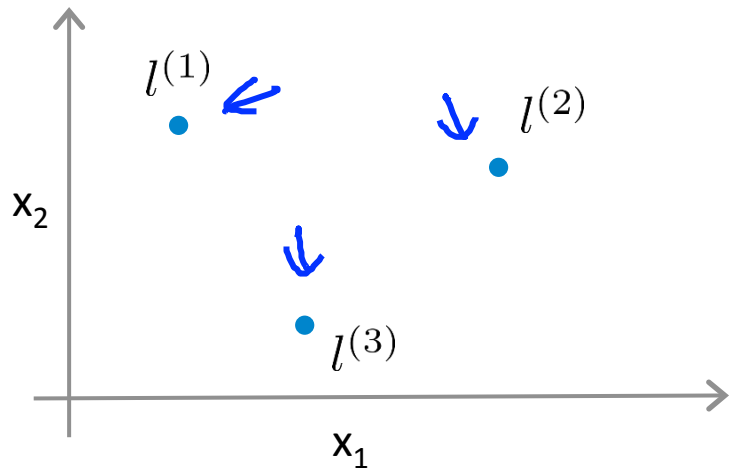
$$h_0(x) = \begin{cases} 1 & \text{if } \theta_0 + \theta_1 x_1 + \dots \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$

$$\rightarrow \theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \dots$$

$$f_1 = x_1, \quad f_2 = x_2, \quad f_3 = x_1 x_2, \quad f_4 = x_1^2, \quad f_5 = x_2^2, \dots$$

Is there a different / better choice of the features f_1, f_2, f_3, \dots ?

Choosing the landmarks

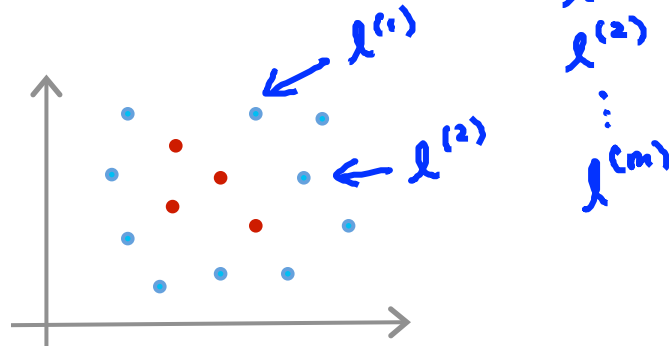
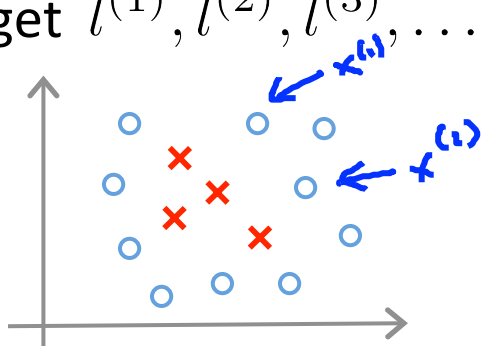


Given x :

$$\begin{aligned} \rightarrow f_i &= \text{similarity}(x, l^{(i)}) \\ &= \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right) \leftarrow \end{aligned}$$

Predict $y = 1$ if $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$ \leftarrow

Where to get $l^{(1)}, l^{(2)}, l^{(3)}, \dots$?



The kernel-based function is exactly equivalent to preprocessing the data by applying similarity function to all inputs, then learning a linear model in the new transformed space.

SVM with Kernels

- Given $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$,
- choose $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(m)} = x^{(m)}$.

Given example x :

$$\begin{aligned} \rightarrow f_1 &= \text{similarity}(x, l^{(1)}) \\ \rightarrow f_2 &= \text{similarity}(x, l^{(2)}) \\ &\dots \end{aligned}$$

$$f = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_n \end{bmatrix} \quad f_0 = 1$$

For training example $(x^{(i)}, y^{(i)})$:

$$\begin{aligned} |x^{(i)} \rightarrow \begin{bmatrix} f_1^{(i)} &= \sin(x^{(i)}, l^{(1)}) \\ f_2^{(i)} &= \sin(x^{(i)}, l^{(2)}) \\ \vdots & \\ f_m^{(i)} &= \sin(x^{(i)}, l^{(m)}) \end{bmatrix} \leftarrow x^{(i)} \right. \\ & \left. f_i^{(i)} = \sin(x^{(i)}, l^{(i)}) = \exp\left(-\frac{0}{2\sigma_i}\right) = 1 \right. \\ & \left. |x^{(i)} \in \mathbb{R}^{n+1} \rightarrow \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_n^{(i)} \end{bmatrix} \text{ (or } \mathbb{R}^n) \right. \\ & \left. f_0^{(i)} = 1 \right. \end{aligned}$$

Commonly used kernels

- Homogeneous polynomials

$$k(x, y) = (\langle x, y \rangle)^d$$

- Inhomogeneous polynomials

$$k(x, y) = (\langle x, y \rangle + 1)^d$$

- Gaussian Kernel

$$k(x, y) = \exp\left(-\frac{\|x - y\|^2}{2\sigma^2}\right)$$

- Sigmoid Kernel

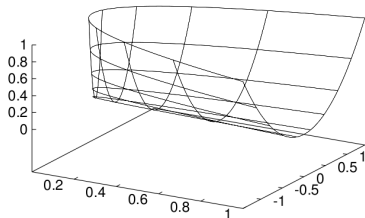
$$k(x, y) = \tanh(\eta \langle x, y \rangle + \nu)$$

Polynomial kernel

$$k(x, y) = (\langle x, y \rangle)^d$$

Example: $n = 2, d = 2, x = (x_1, x_2)$

- $\Phi(x) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$

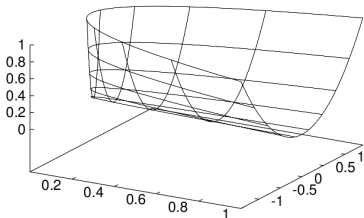


Polynomial kernel

$$k(x, y) = (\langle x, y \rangle)^d$$

Example: $n = 2, d = 2, x = (x_1, x_2)$

- $\Phi(x) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$



- Neither the mapping Φ nor the feature space is unique
 - $\Phi(x) = (x_1^2, x_1x_2, x_1x_2, x_2^2)$
 - $\Phi(x) = \frac{1}{\sqrt{2}} (x_1^2 - x_2^2, 2x_1x_2, x_1^2 + x_2^2)$

Logistic regression vs. SVMs

n = number of features ($x \in \mathbb{R}^{n+1}$), m = number of training examples

→ If n is large (relative to m): (e.g. $n \geq m$, $n = \underline{10,000}$, $m = \underline{10} \dots \underline{1000}$)

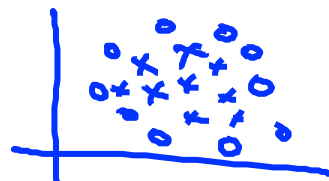
→ Use logistic regression, or SVM without a kernel ("linear kernel")

→ If n is small, m is intermediate: ($n = \underline{1-1000}$, $m = \underline{10-10,000}$) ←

→ Use SVM with Gaussian kernel

If n is small, m is large: ($n = \underline{1-1000}$, $m = \underline{50,000+}$)

→ Create/add more features, then use logistic regression or SVM without a kernel



→ Neural network likely to work well for most of these settings, but may be slower to train.